                  Relayed Echo Reply mechanism for LSP Ping
                   draft-ietf-mpls-lsp-ping-relay-reply-08

Abstract

   In some inter autonomous system (AS) and inter-area deployment
   scenarios for RFC 4379 "Label Switched Path (LSP) Ping and
   Traceroute", a replying LSR may not have the available route to the
   initiator, and the Echo Reply message sent to the initiator would be
   discarded resulting in false negatives or complete failure of
   operation of LSP Ping and Traceroute.  This document describes
   extensions to LSP Ping mechanism to enable the replying Label
   Switching Router (LSR) to have the capability to relay the Echo
   Response by a set of routable intermediate nodes to the initiator.
   This document updates RFC 4379.

Status of This Memo

Copyright Notice

Table of Contents

## 1.  Introduction

   This document describes the extensions to the Label Switched Path
   (LSP) Ping as specified in [RFC4379], by adding a relayed echo reply
   mechanism which could be used to detect data plane failures for the
   inter autonomous system (AS) and inter-area LSPs.  The extensions are
   to update the [RFC4379].  Without these extensions, the ping
   functionality provided by [RFC4379] would fail in many deployed
   inter-AS scenarios, since the replying LSR in one AS may not have the
   available route to the initiator in the other AS.  The mechanism in
   this document defines a new message type referred as "Relayed Echo
   Reply message", and a new TLV referred as "Relay Node Address Stack
   TLV".

   This document is also to update [RFC4379], include updating of Echo
   Request sending procedure in section 4.3 of [RFC4379], Echo Request
   receiving procedure in section 4.4 of [RFC4379], Echo Reply sending
   procedure in Section 4.5 of [RFC4379], Echo Reply receiving procedure
   in section 4.6 of [RFC4379].

## 1.1.  Conventions Used in This Document

   The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT",
   "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this
   document are to be interpreted as described in [RFC2119].

## 2.  Motivation

   LSP Ping [RFC4379] defines a mechanism to detect the data plane
   failures and localize faults.  The mechanism specifies that the Echo
   Reply should be sent back to the initiator using an UDP packet with
   the IPv4/IPv6 address of the originating LSR.  This works in
   administrative domains where IP addresses reachability are allowed
   among LSRs, and every LSR is able to route back to the originating
   LSR.  However, in practice, this is often not the case due to intra-
   provider routing policy, route hiding, and network address
   translation at autonomous system border routers (ASBR).  In fact, it
   is almost uniformly the case that in inter-AS scenarios, it is not

allowed the distribution or direct routing to the IP addresses of any
of the nodes other than the ASBR in another AS.

Figure 1 demonstrates a case where one LSP is set up between PE1 and
PE2.  If PE1's IP address is not distributed to AS2, a traceroute
from PE1 directed towards PE2 can result in a failure because an LSR
in AS2 may not be able to send the Echo Reply message.  E.g., P2
cannot forward packets back to PE1 given that it is an routable IP
address in AS1 but not routable in AS2.  In this case, PE1 would
detect a path break, as the Echo Reply messages would not be
received.  Then localization of the actual fault would not be
possible.

Note that throughout the document, routable address means that it is
possible to route an IP packet to this address using the normal
information exchanged by the IGP operating in the AS.

```
+-------+   +-------+   +------+   +------+   +------+   +------+
|       |   |       |   |      |   |      |   |      |   |      |
|  PE1  +---+   P1  +---+ ASBR1+---+ ASBR2+---+  P2  +---+  PE2 |
|       |   |       |   |      |   |      |   |      |   |      |
+-------+   +-------+   +------+   +------+   +------+   +------+
<--------------AS1------------><--------------AS2----------->
<-------------------------- LSP ---------------------------->
```

Figure 1: Simple Inter-AS LSP Configuration

A second example that illustrates how [RFC4379] would be insufficient
would be the inter-area situation in a seamless MPLS architecture
[I-D.ietf-mpls-seamless-mpls] as shown below in Figure 2.  In this
example LSRs in the core network would not have IP reachable route to
any of the ANs.  When tracing an LSP from one AN to the remote AN,
the LSR1/LSR2 node cannot send the Echo Reply either, like the P2
node in the inter-AS scenario in Figure 1.

```
          +-------+   +-------+   +------+   +------+
          |       |   |       |   |      |   |      |
       +--+ AGN11 +---+ AGN21 +---+ ABR1 +---+ LSR1 +--> to AGN
      /   |       | /|        |   |      |   |      |     |
 +----+/     +-------+\/ +-------+   +------+  /+------+
 | AN |          /\                      \/
 +----+\     +-------+ \+-------+   +------+/\ +------+
      \   |       |   |       |   |      |   | \|     |
       +--+ AGN12 +---+ AGN22 +---+ ABR2 +---+ LSR2 +--> to AGN
          |       |   |       |   |      |   |      |     |
          +-------+   +-------+   +------+   +------+
   static route    ISIS L1 LDP            ISIS L2 LDP
   <-Access-><--Aggregation Domain--><---------Core--------->
```

Figure 2: Seamless MPLS Architecture

This document describes extensions to the LSP Ping mechanism to
facilitate a response from the replying LSR, by defining a mechanism
that uses a relay node (e.g, ASBR) to relay the message back to the
initiator.  Every designated or learned relay node must be reachable
to the next relay node or to the initiator.  Using a recursive
approach, relay node could relay the message to the next relay node
until the initiator is reached.

The LSP Ping relay mechanism in this document is defined for unicast
case.  How to apply the LSP Ping relay mechanism in multicast case is
out of the scope.

## 3.  Extensions

[RFC4379] describes the basic MPLS LSP Ping mechanism, which defines
two message types, Echo Request and Echo Reply message.  This
document defines a new message, Relayed Echo Reply message.  This new
message is used to replace Echo Reply message which is sent from the
replying LSR to a relay node or from a relay node to another relay
node.

A new TLV named Relay Node Address Stack TLV is defined in this
document, to carry the IP addresses of the possible relay nodes for
the replying LSR.

In addition, MTU (Maximum Transmission Unit) Exceeded Return Code is
defined to indicate to the initiator that one or more TLVs will not
be returned due to MTU size.

It should be noted that this document focuses only on detecting the
LSP which is set up using a uniform IP address family type.  That is,

all hops between the source and destination node use the same address
family type for their LSP ping control planes.  This does not
preclude nodes that support both IPv6 and IPv4 addresses
simultaneously, but the entire path must be addressable using only
one address family type.  Supporting for mixed IPv4-only and
IPv6-only is beyond the scope of this document.

## 3.1.  Relayed Echo Reply message

The Relayed Echo Reply message is a UDP packet, and the UDP payload
has the same format with Echo Request/Reply message.  A new message
type is requested from IANA.

```
New Message Type:
    Value    Meaning
    -----    -------
    TBD      MPLS Relayed Echo Reply
```

The use of TCP and UDP port number 3503 is described in [RFC4379] and
has been allocated by IANA for LSP Ping messages.  The Relayed Echo
Reply message will use the same port number.

## 3.2.  Relay Node Address Stack

The Relay Node Address Stack TLV is an optional TLV.  It MUST be
carried in the Echo Request, Echo Reply and Relayed Echo Reply
messages if the echo reply relayed mechanism described in this
document is required.  Figure 3 illustrates the TLV format.

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|              Type             |             Length            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|    Initiator Source Port      | Reply Add Type|   Reserved    |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Source Address of Replying Router (0, 4, or 16 octects)   |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Destination Address Pointer   |   Number of Relayed Addresses |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                                                               |
~              Stack of Relayed Addresses                       ~
|                                                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

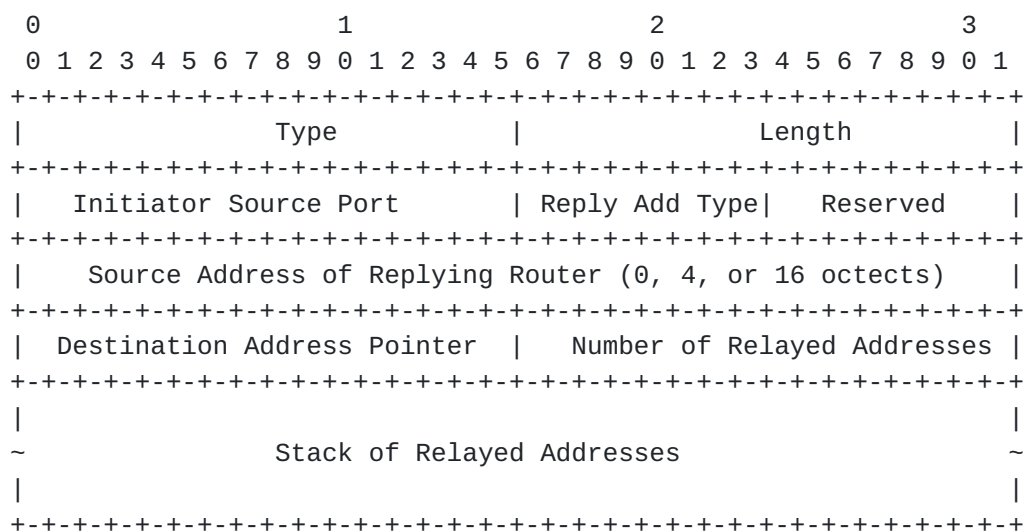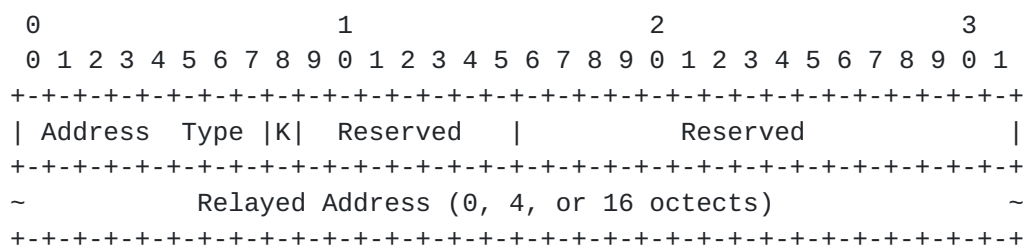Figure 3: Relay Node Address Stack TLV

- Type: to be assigned by IANA.  A value should be assigned from
  32768-49161 as suggested by [RFC4379] Section 3.

- Length: the length of the value field in octets.

- Initiator Source Port: the source UDP port that the initiator uses
  in the Echo Request message, and also the port that is expected to
  receive the Echo Reply message.

- Reply Address Type: address type of replying router.  This value
  also implies the length of the address field as shown below.

```
Type#   Address Type   Address Length
----    ------------   ------------
0       Null           0
1       IPv4           4
2       IPv6           16
```

- Reserved: This field is reserved and MUST be set to zero.

- Source Address of Replying Router: source IP address of the
  originator of Echo Reply or Replay Echo Reply message.

- Destination Address Pointer: an integer entry number used as the
  destination address of the Reply or Relayed Reply message.  The
  entry on the top of the Stack of Relayed Addresses will have value
  1.

- Number of Relayed Addresses: an integer indicating the number of
  relayed addresses in the stack.

- Stack of Relayed Addresses: a list of relay node addresses.

The format of each relay node address is as below:

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Address  Type |K|  Reserved   |        Reserved              |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
~           Relayed Address (0, 4, or 16 octects)              ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

```
Type#   Address Type   Address Length
----    ------------   ------------
0       Null           0
1       IPv4           4
2       IPv6           16
```

Reserved: The two fields are reserved and MUST be set to zero.

K bit: if the K bit is set to 1, then this sub-TLV MUST be kept in Relay Node Address Stack during TLV compress process described in [section 4.2](#).

Having the K bit set in the relay node address entry causes that entry to be preserved in the Relay Node Address Stack TLV for the entire traceroute operation.  A responder node MAY set the K bit to ensure its relay node address entry remains as one of the relay nodes in the Relay Node Address Stack TLV.  The address with K bit set will always be a relay node address for the Relayed Echo Reply, see [section 4.3](#).

Relayed Address: this field specifies the node address, either IPv4 or IPv6.

### 3.3.  MTU Exceeded Return Code

Return Code is defined to indicate that one or more TLVs were omitted from the Echo Reply or Relayed Echo Reply message to avoid exceeding the message's effective MTU size.  These TLVs MAY be included in an Errored TLV's Object with their lengths set to 0 and no value.  The return sub-code MUST be set to the value that otherwise would have been sent.

```
MTU Exceeded Return Code:
    Value    Meaning
    -----    -------
    TBD      One or more TLVs not returned due to MTU size
```

### 4.  Procedures

### 4.1.  Sending an Echo Request

In addition to the procedures described in [section 4.3 of [RFC4379]](#), a Relay Node Address Stack TLV MUST be carried in the Echo Request message to facilitate the relay functionality.

When the initiator sends the first Echo Request with a Relay Node Address Stack TLV, the TLV MUST contain the initiator address as the first entry of the stack of relayed addresses, the destination address pointer set to this entry, and the source address of the replying router set to null.  The source UDP port field MUST be set to the source UDP port.  Note that the first relay node address in the stack will always be the initiator's address.

4.2.  Receiving an Echo Request

   An LSR that does not recognize the Relay Node Address Stack TLV,
   SHOULD ignore it as per section 3 of [RFC4379].

   In addition to the processes in section 4.4 of [RFC4379], the
   procedures of the Relay Node Address Stack TLV are defined here.

   Upon receiving a Relay Node Address Stack TLV in an Echo Request
   message, the receiver updates the "Source Address of Replying
   Router".  The address MUST be same as the source IP address of Relay
   Echo Reply (section 4.3) or Echo Reply message (section 4.5) being
   sent.

   Those address entries with K bit set to 1 MUST be kept in the stack.
   The receiver MUST check the addresses of the stack in sequence from
   bottom to top to find the last address in the stack with the K bit
   set (or the top of the stack if no K bit was found).  The receiver
   then checks the stack beginning with this entry, proceeding towards
   the bottom to find the first routable address IP address.  The
   Destination Address Pointer MUST be set to this entry which is also
   the Destination Address.  Address entries below the first routable IP
   address MUST be deleted.  At least one address entries of the
   replying LSR MUST be added at the bottom of the stack.  A second or
   more address entries MAY also be added if necessary, depending on
   implementation.  The final address added MUST be an address that is
   reachable through the interface that the Echo Request Message would
   have been forwarded if it had not TTL expired at this node.  The
   updated Relay Node Address Stack TLV MUST be carried in the response
   message.

   If the replying LSR is configured to hide its routable address
   information, the address entry added in the stack MUST be a NIL entry
   with Address Type set to NULL.

   If a node spans two addressing domains (with respect to this message)
   where nodes on either side may not be able to reach nodes in the
   other domain, then the final address added SHOULD set the K bit.  K
   bit applies in the case of a NULL address, to serve as a warning to
   the initiator that further Echo Request messages may not result in
   receiving Echo Reply messages.

   If the full reply message would exceed the MTU size, the Relay Node
   Address Stack TLV SHOULD be included in the Echo Reply message (i.e.,
   other optional TLVs are excluded).

### 4.3.  Originating an Relayed Echo Reply

The Destination Address determined in section 4.2 is used as the next
relay node address.  If the resolved next relay node address is not
routable, then sending of Relayed Echo Reply or Echo Reply will fail.

If the first IP address in the Relay Node Address Stack TLV is not
the next relay node address, the replying LSR SHOULD send a Relayed
Echo Reply message to the next relay node.  The processing of Relayed
Echo Reply is the same with the procedure of the Echo Reply described
in Section 4.5 of [RFC4379], except the destination IP address and
the destination UDP port.  The destination IP address of the Relayed
Echo Reply is set to the next relay node address from the Relay Node
Address Stack TLV, and both the source and destination UDP port are
set to 3503.  The updated Relay Node Address Stack TLV described in
section 4.2 MUST be carried in the Relayed Echo Reply message.  The
Source Address of Replying Router field is kept unmodified, and
Source IP address field of the IP header is set to an address of the
sending node.

### 4.4.  Relaying an Relayed Echo Reply

Upon receiving an Relayed Echo Reply message with its own address as
the destination address in the IP header, the relay node MUST
determine the next relay node address as described in section 4.2,
with the modification that the location of the received Destination
Address is used instead of the bottom of stack in the algorithm.  The
Destination Address Pointer in Relay Node Address Stack TLV will be
set to the next relay node address.  Note that unlike section 4.2 no
changes are made to the Stack of Relayed Addresses.

If the next relay node address is not the first one in the address
list, i.e., another intermediate relay node, the relay node MUST send
an Relayed Echo Reply message to the determined upstream node with
the payload unchanged other than the Relay Node Address Stack TLV.
The TTL SHOULD be copied from the received Relay Echo Reply and
decremented by 1.  The Source Address of Replying Router field is
kept unmodified, and Source IP address field of the IP header is set
to an address of the sending node.

### 4.5.  Sending an Echo Reply

The Echo Reply is sent in two cases:

1.  When the replying LSR receives an Echo Request, and the first IP
address in the Relay Node Address Stack TLV is the next relay node
address (section 4.3), the replying LSR would send an Echo Reply to
the initiator.  In addition to the procedure of the Echo Reply

described in Section 4.5 of [RFC4379], the updated Relay Node Address
Stack TLV described in section 4.2 MUST be carried in the Echo Reply.

2.  When the intermediate relay node receives a Relayed Echo Reply,
and the first IP address in the Relay Node Address Stack TLV is the
next relay node address (section 4.4), the intermediate relay node
would send the Echo Reply to the initiator, and update the Message
Type field from type of Relayed Echo Reply to Echo Reply.  The
updated Relay Node Address Stack TLV described in section 4.4 MUST be
carried in the Echo Reply.  The destination IP address of the Echo
Reply is set to the first IP address in the stack, and the
destination UDP port would be copied from the Initiator Source Port
field of the Relay Node Address Stack TLV.  The source UDP port
should be 3503.  The TTL of the Echo Reply SHOULD be copied from the
received Relay Echo Reply and decremented by 1.  The Source Address
of Replying Router field is kept unmodified, and Source IP address
field of the IP header is set to an address of the sending node.

## 4.6.  Sending Subsequent Echo Requests

During a traceroute operation, multiple Echo Request messages are
sent.  Each time the TTL is increased, the initiator could copy the
Relay Node Address Stack TLV received in the previous Echo Reply to
the next Echo Request.Some modifications could also be made to the
stack TLV.  The NIL entry with K bit set SHOULD be deleted, otherwise
the Echo Reply message will not be returned.  The fields of Source
Address of Replying Router and Destination Address Pointer may be
preserved or may be reset for subsequent MPLS Echo Request, and to be
ignored in received MPLS Echo Request.

## 4.7.  Impact to Traceroute

Source IP address in Echo Reply and Relay Echo Reply is to be of the
address of the node sending those packets, not the original
responding node.  Then the traceroute address output module will
print the source IP address as below:

```
if (Relay Node Address Stack TLV exists) {
      Print the Source Address of Replying Router in
      Relay Node Address Stack TLV;
} else {
      Print the source IP address of Echo Reply message;
}
```

5. **LSP Ping Relayed Echo Reply Example**

   Considering the inter-AS scenario in Figure 4 below.  AS1 and AS2 are
   two independent address domains.  In the example, an LSP has been
   created between PE1 to PE2, but PE1 in AS1 is not reachable by P2 in
   AS2.


```
+-------+   +-------+   +------+   +------+   +------+   +------+
|       |   |       |   |      |   |      |   |      |   |      |
| PE1  +---+  P1  +---+ ASBR1+---+ ASBR2+---+  P2  +---+  PE2 |
|       |   |       |   |      |   |      |   |      |   |      |
+-------+   +-------+   +------+   +------+   +------+   +------+
<---------------AS1------------><---------------AS2------------>
<------------------------ LSP ------------------------------->
```
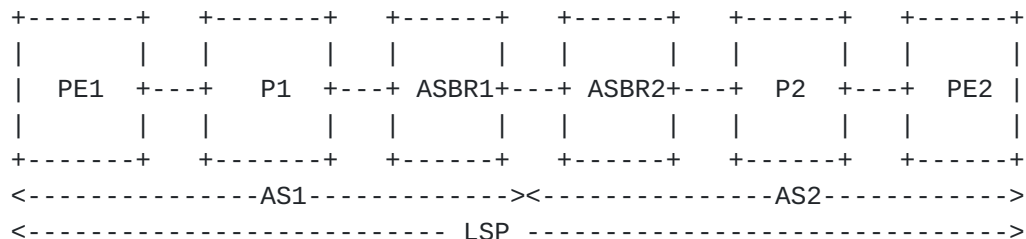

                   Figure 4: Example Inter-AS LSP

   When performing LSP traceroute on the LSP, the first Echo Request
   sent by PE1 with outer-most label TTL=1, contains the Relay Node
   Address Stack TLV with PE1's address as the first relayed address.

   After processed by P1, P1's interface address facing ASBR1 without
   the K bit set will be added in the Relay Node Address Stack TLV
   address list following PE1's address in the Echo Reply.

   PE1 copies the Relay Node Address Stack TLV into the next Echo
   Request when receiving the Echo Reply.

   Upon receiving the Echo Request, ASBR1 checks the address list in the
   Relay Node Address Stack TLV, and determines that PE1's address is
   the next relay address.  Then deletes P1's address, and adds its
   interface address facing ASBR2 with the K bit set.  As a result,
   there would be PE1's address followed by ASBR1's interface address
   facing ASBR2 in the Relay Node Address Stack TLV of the Echo Reply
   sent by ASBR1.

   PE1 then sends an Echo Request with outer-most label TTL=3,
   containing the Relay Node Address Stack TLV copied from the received
   Echo Reply message.  Upon receiving the Echo Request message, ASBR2
   checks the address list in the Relay Node Address Stack TLV, and
   determines ASBR1's interface address is the next relay address in the
   stack TLV.  ASBR2 adds its interface address facing P2 with the K bit
   set.  Then ASBR2 sets the next relay address as the destination
   address of the Relayed Echo Reply, and sends the Relayed Echo Reply
   to ASBR1.

Upon receiving the Relayed Echo Reply from ASBR2, ASBR1 checks the
address list in the Relay Node Address Stack TLV, and determines that
PE1's address is the next relay node.  Then ASBR1 sends an Echo Reply
to PE1.

For the Echo Request with outer-most label TTL=4, P2 checks the
address list in the Relay Node Address Stack TLV, and determines that
ASBR2's interface address is the next relay address.  Then P2 sends
an Relayed Echo Reply to ASBR2 with the Relay Node Address Stack TLV
containing four addresses, PE1's, ASBR1's interface address, ASBR2's
interface address and P2's interface address facing PE2 in sequence.

Then according to the process described in section 4.4, ASBR2 sends
the Relayed Echo Reply to ASBR1.  Upon receiving the Relayed Echo
Reply, ASBR1 sends an Echo Reply to PE1.  And as relayed by ASBR2 and
ASBR1, the Echo Reply would finally be sent to the initiator PE1.

For the Echo Request with outer-most label TTL=5, the Echo Reply
would relayed to PE1 by ASBR2 and ASBR1, similar to the case of
TTL=4.

The Echo Reply from the replying node which has no IP reachable route
to the initiator is thus transmitted to the initiator by multiple
relay nodes.

## 6.  Security Considerations

The Relayed Echo Reply mechanism for LSP Ping creates an increased
risk of DoS by putting the IP address of a target router in the Relay
Node Address Stack.  These messages then could be used to attack the
control plane of an LSR by overwhelming it with these packets.  A
rate limiter SHOULD be applied to the well-known UDP port on the
relay node as suggested in [RFC4379].  The node which acts as a relay
node SHOULD validate the relay reply against a set of valid source
addresses and discard packets from untrusted border router addresses.
An implementation SHOULD provide such filtering capabilities.

If an operator wants to obscure their nodes, it is RECOMMENDED that
they may replace the replying node address that originated the Echo
Reply with NIL address entry in Relay Node Address Stack TLV.

A receiver of an MPLS Echo Request could verify that the first
address in the Relay Node Address Stack TLV is the same address as
the source IP address field of the received IP header.

Other security considerations discussed in [RFC4379], are also
applicable to this document.

7.  Backward Compatibility

   When one of the nodes along the LSP does not support the mechanism
   specified in this document, the node will ignore the Relay Node
   Address Stack TLV as described in section 4.2.  Then the initiator
   may not receive the Relay Node Address Stack TLV in Echo Reply
   message from that node.  In this case, an indication should be
   reported to the operator, and the Relay Node Address Stack TLV in the
   next Echo Request message should be copied from the previous Echo
   Request, and continue the ping process.  If the node described above
   is located between the initiator and the first relay node, the ping
   process could continue without interruption.

8.  IANA Considerations

   IANA is requested to assign one new Message Type, one new TLV and one
   Return Code.

8.1.  New Message Type

   This document requires allocation of one new message type from
   "Multi-Protocol Label Switching (MPLS) Label Switched Paths (LSPs)
   Ping Parameters" registry, the "Message Type" registry:

        Value    Meaning
        -----    -------
        TBD      MPLS Relayed Echo Reply

   The value should be assigned from the "Standards Action" range
   (0-191), and using the lowest free value within this range.

8.2.  New TLV

   This document requires allocation of one new TLV from "Multi-Protocol
   Label Switching (MPLS) Label Switched Paths (LSPs) Ping Parameters"
   registry, the "TLVs" registry:

        Type     Meaning
        ----     --------
        TBD      Relay Node Address Stack TLV

   A suggested value should be assigned from "Standards Action" range
   (32768-49161) as suggested by [RFC4379] Section 3, using the first
   free value within this range.

### 8.3.  MTU Exceeded Return Code

   This document requires allocation of MTU Exceeded return code from
   "Multi-Protocol Label Switching (MPLS) Label Switched Paths (LSPs)
   Ping Parameters" registry, the "Return Codes" registry:

```
   Value    Meaning
   -----    -------
   TBD      One or more TLVs not returned due to MTU size
```

   The value should be assigned from the "Standards Action" range
   (0-191), and using the lowest free value within this range.

### 9.  Acknowledgement

   The authors would like to thank Carlos Pignataro, Xinwen Jiao, Manuel
   Paul, Loa Andersson, Wim Henderickx, Mach Chen, Thomas Morin, Gregory
   Mirsky, Nobo Akiya and Joel M.  Halpern for their valuable comments
   and suggestions.

### 10.  Contributors

```
   Ryan Zheng
   JSPTPD
   371, Zhongshan South Road
   Nanjing, 210006, China
   Email: ryan.zhi.zheng@gmail.com
```

### 11.  References

### 11.1.  Normative References

   [RFC2119]  Bradner, S., "Key words for use in RFCs to Indicate
              Requirement Levels", BCP 14, RFC 2119, March 1997.

   [RFC4379]  Kompella, K. and G. Swallow, "Detecting Multi-Protocol
              Label Switched (MPLS) Data Plane Failures", RFC 4379,
              February 2006.

### 11.2.  Informative References

   [I-D.ietf-mpls-seamless-mpls]
              Leymann, N., Decraene, B., Filsfils, C., Konstantynowicz,
              M., and D. Steinberg, "Seamless MPLS Architecture", draft-
              ietf-mpls-seamless-mpls-07 (work in progress), June 2014.

Authors' Addresses

    Jian Luo (editor)
    ZTE
    50, Ruanjian Avenue
    Nanjing, 210012, China

    Email: luo.jian@zte.com.cn


    Lizhong Jin (editor)
    Individual
    Shanghai, China

    Email: lizho.jin@gmail.com


    Thomas Nadeau (editor)
    Lucidvision

    Email: tnadeau@lucidvision.com


    George Swallow (editor)
    Cisco
    300 Beaver Brook Road
    Boxborough , MASSACHUSETTS 01719, USA

    Email: swallow@cisco.com