

Network Working Group
Internet-Draft
Expires: June 24, 2006

J. Arkko
Ericsson
I. Beijnum
Muada
December 21, 2005

Failure Detection and Locator Pair Exploration Protocol for IPv6
Multihoming
draft-ietf-shim6-failure-detection-03

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with Section 6 of 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/1id-abstracts.txt>.

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

This Internet-Draft will expire on June 24, 2006.

Copyright Notice

Copyright (C) The Internet Society (2005).

Abstract

This document defines a mechanism for the detection of communication failures between two communicating hosts at IP layer, and an exploration protocol for switching to another pair of interfaces and/or addresses between the same hosts if a working pair can be found. The draft also discusses the roles of a multihoming protocol versus network attachment functions at IP and link layers.

Table of Contents

1. Introduction	3
2. Requirements language	4
3. Related Work	5
4. Definitions	7
4.1. Available Addresses	7
4.2. Locally Operational Addresses	7
4.3. Operational Address Pairs	8
4.4. Current Address Pair	9
4.5. Miscellaneous	10
5. Protocol Overview	11
5.1. Failure Detection	11
5.2. Alternative Address Pair Exploration	13
5.3. Exploration Order	14
5.4. Protocol Design	14
5.5. Example Protocol Runs	15
5.6. Limitations	21
6. Protocol Definition	22
6.1. Keepalive Message	22
6.1.1. Keepalive Option	23
6.2. Probe Message	24
6.2.1. Probe Option	25
6.3. Reachability Option	26
6.3.1. Payload Reception Report	27
6.3.2. Probe Reception Report	28
6.4. Behaviour	30
6.5. Protocol Constants	34
7. Security Considerations	35
8. IANA Considerations	37
9. References	38
9.1. Normative References	38
9.2. Informative References	38
Appendix A. Contributors	40
Appendix B. Acknowledgements	41
Authors' Addresses	42
Intellectual Property and Copyright Statements	43

1. Introduction

The SHIM6 protocol extends IPv6 to support multihoming. This protocol is an IP layer mechanism that hides multihoming from applications [18]. A part of the SHIM6 solution involves detecting when a currently used pair of addresses (or interfaces) between two communication hosts has failed, and picking another pair when this occurs. We call the former failure detection, and the latter locator pair exploration.

This draft defines the mechanism and protocol to achieve both failure detection and locator pair exploration. This protocol is called REACHability Protocol (REAP). It designed to be carried within the SHIM6 protocol, but may also be used in other contexts.

The draft is structured as follows: Section 3 discusses prior work in this space, Section 4 defines a set of useful terms, Section 5 gives an overview of REAP, and Section 6 specifies the message formats and behaviour in detail. Section 7 discusses the security considerations of REAP.

For the purposes of this draft, we consider an address to be synonymous with a locator. We assume that there are other, higher level identifiers such as CGA public keys or HBA bindings that tie the different locators used by a node together [17].

2. Requirements language

In this document, the key words "MAY", "MUST", "MUST NOT", "OPTIONAL", "RECOMMENDED", "SHOULD", and "SHOULD NOT", are to be interpreted as described in [2].

3. Related Work

In SCTP [10], the addresses of the endpoints are learned in the connection setup phase either through listing them explicitly or via giving a DNS name that points to them. In order to provide a failover mechanism between multihomed hosts, SCTP selects one of the peer's addresses as the primary address by the application running on top of SCTP. All data packets are sent to this address until there is a reason to choose another address, such as the failure of the primary address.

SCTP also tests the reachability of the peer endpoint's addresses. This is done both via observing the data packets sent to the peer or via a periodic heartbeat when there is no data packets to send. Each time data packet retransmission is initiated (or when a heartbeat is not answered within the estimated round-trip time) an error counter is incremented. When a configured error limit is reached, the particular destination address is marked as inactive. The reception of an acknowledgement or heartbeat response clears the counter. Retransmission: When retransmitting the endpoint attempts pick the most "divergent" source-destination pair from the original source- destination pair to which the packet was transmitted. Rules for such selection are, however, left as implementation decisions in SCTP.

SCTP does not define how local knowledge (such as information learned from the link layer) should be used. SCTP also has no mechanism to deal with dynamic changes to the set of available addresses, although mechanisms for that are being developed [20].

The MOBIKE protocol [15] provides multihoming and mobility for VPN connections. Its failure detection and locator pair exploration is designed to work across mixed IPv4/IPv6 environments and NATs, as long as a path that allows bidirectional communication can be found.

Existing mechanisms at lower layers or in IKEv2 are used to detect failures, and upon failure MOBIKE attempts to explore all combinations of addresses to find a working pair. Such exploration is necessary when a problem affects both nodes. For instance, two nodes connected by two separate point-to-point links will be unable to switch to the other link if a failure occurs on the first one. While both communicating hosts are aware of each others' addresses, only one end of the communication is in charge of deciding what address pair to use, however.

The mobility and multihoming specification for the HIP protocol [14] leaves the determination of when address updates are sent to a local policy, but suggests the use of local information and ICMP error messages.

Network attachment procedures are also relevant for multihoming. The IPv6 and MIPv6 working groups have standardized mechanisms to learn about networks that a node has attached to. Basic IPv6 Neighbor Discovery was, however, designed primarily for static situations. The fully dynamic detection procedure has turned out to be a relatively complex procedure for mobile hosts, and it was not fully anticipated at the time IPv6 Neighbor Discovery or DHCP were being designed. As a result, enhanced or optimized mechanisms are being designed in the DHC and DNA working groups [13] [7].

ICE [16], STUN [11], and TURN [24] are also related mechanisms. They are primarily used for NAT detection and communication through NATs in IPv4 environment, for application such as voice over IP. STUN uses a server in the Internet to discover the presence and type of NATs and the client's public IP addresses and ports. TURN makes it possible to receive incoming connections in hosts behind NATs. ICE makes use of these protocols in peer-to-peer cooperative fashion, allowing participants to discover, create and verify mutual connectivity, and then use this connectivity for multimedia streams. While these mechanisms are not designed for dynamic and failure situations, they have many of the same requirements for the exploration of connectivity, as well as the requirement to deal with middleboxes.

Related work in the IPv6 area includes RFC 3484 [6] which defines source and destination address selection rules for IPv6 in situations where multiple candidate address pairs exist. RFC 3484 considers only a static situation, however, and does not take into account the effect of failures. Reference [23] considers how applications can re-initiate connections after failures in the best way. This work differs from the shim-layer approach selected for further development in the working group with respect to the timing of the address selection. In the shim-layer approach failure detection and the selection of new addresses happens at any time, while [23] considers only the case when an application re-establishes connections.

An earlier SHIM6 document [19] discussed what kind of mechanisms can be used to detect whether the peer is still reachable at the currently used address. Two proposed mechanisms, Correspondent Unreachability Detection (CUD) and Forced Bidirectional Communication (FBD) were presented. CUD is based on getting upper layer positive feedback, and IPv6 NUD-like probing if there is no feedback. FBD is based on forcing bidirectional communication by adding keepalive messages when there is no other, payload traffic. FBD is the chosen mechanism in this document.

4. Definitions

This section defines terms useful in discussing the problem space.

4.1. Available Addresses

Multihoming nodes need to be aware of what addresses they themselves have. If a node loses the address it is currently using for communications, another address must replace this address. And if a node loses an address that the node's peer knows about, the peer must be informed. Similarly, when a node acquires a new address it may generally wish the peer to know about it.

Definition. Available address. An address is said to be available if the following conditions are fulfilled:

- o The address has been assigned to an interface of the node.
- o If the address is an IPv6 address, we additionally require that (a) the address is valid in the sense of RFC 2461 [3], and that (b) the address is not tentative in the sense of RFC 2462 [4]. In other words, the address assignment is complete so that communications can be started.

Note that this explicitly allows an address to be optimistic in the sense of [8] even though implementations are probably better off using other addresses as long as there is an alternative.

- o The address is a global unicast, unique local address [9], or an unambiguous IPv6 link-local address. That is, it is not an IPv6 site-local address. Where IPv6 link-local addresses are used, their use needs to be unambiguous as follows. At most one link- local address may be used per node within the same connection between two peers.
- o The address and interface is acceptable for use according to a local policy.

Available addresses are discovered and monitored through mechanisms outside the scope of the protocol described here. These mechanisms include IPv6 Neighbor Discovery and Address Autoconfiguration [3] [4], DHCP [5], and DNS mechanisms [7].

4.2. Locally Operational Addresses

Two different granularity levels are needed for failure detection. The coarser granularity is for individual addresses:

Definition. **Locally Operational Address.** An available address is said to be locally operational when its use is known to be possible locally: the interface is up, at least one default router (if applicable) that could be used to send a packet with this address as a source address is known to be reachable, and no other local information points to the address being unusable.

Locally operational addresses are discovered and monitored through mechanisms outside the protocol described here. These mechanisms include IPv6 Neighbor Discovery [3] and link layer specific mechanisms.

It is also possible for hosts to learn about routing failures for a particular selected source prefix, if suitable protocols for this purpose exist. Some proposals in this space have been made, see, for instance [21] and [23]. Potential approaches include overloading information in current IPv6 Router Advertisement or adding some new information in them. Similarly, hosts could learn information from servers that query the BGP routing tables.

4.3. Operational Address Pairs

The existence of locally operational addresses are not, however, a guarantee that communications can be established with the peer. A failure in the routing infrastructure can prevent the sent packets from reaching their destination. For this reason we need the definition of a second level of granularity, for pairs of addresses:

Definition. **Bidirectionally operational address pair.** A pair of locally operational addresses are said to be an operational address pair, iff bidirectional connectivity can be shown between the addresses. That is, a packet sent with one of the addresses in the source field and the other in the destination field reaches the destination, and vice versa.

Unfortunately, there are scenarios where bidirectionally operational address pairs do not exist. For instance, ingress filtering or network failures may result in one address pair being operational in one direction while another one is operational from the other direction. The following definition captures this general situation:

Definition. **Unidirectionally operational address pair.** A pair of locally operational addresses are said to be an unidirectionally operational address pair, iff packets sent with the first address as the source and the second address as the destination can be shown to reach the destination.

Both types of operational pairs could be discovered and monitored

through the following mechanisms:

- o Positive feedback from upper layer protocols. For instance, TCP can indicate to the IP layer that it is making progress. This is similar to how IPv6 Neighbor Unreachability Detection can in some cases be avoided when upper layers provide information about bidirectional connectivity [3]. In the case of unidirectional connectivity, the upper layer protocol responses come back using another address pair, but show that the messages sent using the first address pair have been received.
- o Negative feedback from upper layer protocols. It is conceivable that upper layer protocols give an indication of a problem to the multihoming layer. For instance, TCP could indicate that there's either congestion or lack of connectivity in the path because it is not getting ACKs.
- o Explicit reachability tests, such as keepalives or probes added when there's only unidirectional payload traffic.
- o ICMP error messages. Given the ease of spoofing ICMP messages, one should be careful to not trust these blindly, however. Our suggestion is to use ICMP error messages only as a hint to perform an explicit reachability test, but not as a reason to disrupt ongoing communications without other indications of problems. The situation may be different when certain verifications of the ICMP messages are being performed [22]. These verifications can ensure that (practically) only on-path attackers can spoof the messages.

Note a multihoming protocol needs to perform a return routability test of an address before it is taken into use. The purpose of this test is to ensure that fraudulent peers do not trick others into redirecting traffic streams onto innocent victims [25]. This test can at the same time work as a means to ensure that an address pair is operational, as discussed in Section 5.2.

4.4. Current Address Pair

IP-layer solutions need to avoid sending packets concurrently over multiple paths; TCP behaves rather poorly in such circumstances. For this reason it is necessary to choose a particular pair of addresses as the current address pair which is used until problems occur, at least for the same session.

A current address pair need not be operational at all times. If there is no traffic to send, we may not know if the primary address pair is operational. Nevertheless, it makes sense to assume that the address pair that worked in some time ago continues to work for new

communications as well.

4.5. Miscellaneous

Addresses can become deprecated [3]. When other operational addresses exist, nodes generally wish to move their communications away from the deprecated addresses.

Similarly, IPv6 source address selection [6] may guide the selection of a particular source address - destination address pair.

5. Protocol Overview

This section discusses the design of the reachability detection and address pair exploration mechanisms, and gives an overview of the REAP protocol.

A naive implementation of an (un)reachability detection mechanism could just probe all possible paths between two hosts periodically. A "path" is defined as a combination of a source address for host A and a destination address for host B. In hop-by-hop forwarding the source address has no effect on reachability, but in the presence of filters or source address based routing, it may. And although links almost always work in two directions, routing protocols and filters only work in one direction so unidirectional reachability is possible. Without additional mechanisms, the practice of ingress filtering by ISPs makes unidirectional connectivity likely. Being able to use the working leg in a unidirectional path is useful, it's not an essential requirement. It is essential, however, to avoid assuming bidirectional connectivity when there is in fact a unidirectional failure.

Exploring the full set of communication options between two hosts that both have two or more addresses is an expensive operation as the number of combinations to be explored increases very quickly with the number of addresses. For instance, with two addresses on both sides, there are four possible address pairs. Since we can't assume that reachability in one direction automatically means reachability for the complement pair in the other direction, the total number of two-way combinations is eight. (Combinations = $n_A * n_B * 2$.)

An important observation in multihoming is that failures are relatively infrequent, so that a path that worked a few seconds ago is very likely to work now as well. So it makes sense to have a light-weight protocol that confirms existing reachability, and only invoke heavier exploration when there is a suspected failure.

5.1. Failure Detection

This process consists of three tasks. First, it is necessary to track local information from lower and upper layers. For instance, when link layer informs that we have no connection then we know there is a failure. Nodes SHOULD employ techniques listed in Section 4.1 and Section 4.2 to be aware of the local situation.

Similarly, it is necessary to track remote address information from the peer. For instance, the peer may inform that its currently used address is no longer in use. Techniques outside the scope of this document are used for this, for further information see [18].

The third task is to ensure verify reachability with the peer when the local and remote information indicates that communication should be possible. This needs to be performed only if there's upper layer packets to be sent, however.

This document defines the protocol mechanisms only for the third task. We employ a technique called Forced Bidirectional Detection (FBD). Reachability for the currently used address pair in a shim context is determined by making sure that whenever there is data traffic in one direction, there is also traffic in the other direction. This can be data traffic as well, but also transport layer acknowledgments or a REAP reachability keepalive if there is no other traffic. This way, it is no longer possible to have traffic in only one direction, so whenever there is data traffic going out, but there are no return packets, there must be a failure, so the full path exploration mechanism is started.

A more detailed description of the current pair reachability evaluation mechanism:

1. The base timing unit for this mechanism is named Keepalive Timeout. Until a negotiation mechanism to negotiate different values for this timer becomes available, the value (3 seconds) specified in Section 6.5 SHOULD be used.
2. Whenever outgoing data packets are generated that are part of a shim context, a timer is started to reflect the requirement that the peer should generate return traffic from data packets.
3. Whenever incoming data packets are received that are part of a shim context, the timer associated with the return traffic from the peer is stopped, and another timer is started to reflect the requirement for this node to generate return traffic.
4. The reception of a REAP keepalive packet leads to stopping the timer associated with the return traffic from the peer.
5. Keepalive Timeout seconds after the last data packet has been received for a context, and if no other packet has been sent within this context since the data packet has been received, a REAP keepalive packet is generated for the context in question and transmitted to the correspondent. A host may send the keepalive sooner than Keepalive Timeout seconds if implementation considerations warrant this. The average time after which keepalives are sent MUST be at least $\text{Keepalive Timeout} / 2$ seconds. After sending a single keepalive message, no additional keepalive messages are sent until a data packet is received within this shim context. Keepalives are not sent at all when a

data packet was sent since the last received data packet.

6. Send Timeout seconds (10 s; see Section 6.5) after the transmission of a data packet with no return traffic on this context, a full reachability exploration is started. This timeout period is larger than the Keepalive Timeout to accommodate for lost keepalives and regular variations in round trip times.

5.2. Alternative Address Pair Exploration

As explained in previous section, the currently used address pair may become invalid either through one of the addresses being becoming unavailable or inoperational, or the pair itself being declared inoperational. An exploration process attempts to find another operational pair so that communications can resume.

What makes this process hard is the requirement to support unidirectionally operational address pairs. It is insufficient to probe address pairs by a simple request - response protocol. Instead, the party that first detects the problem starts a process where it tries each of the different address pairs in turn by sending a message to its peer. These messages carry information about the state of connectivity between the peers, such as whether the sender has seen any traffic from the peer recently. When the peer receives a message that indicates a problem, it assists the process by starting its own parallel exploration to the other direction, again sending information about the recently received payload traffic or signaling messages.

Specifically, when A decides that it needs to explore for an alternative address pair to B, it will initiate a set of Probe messages, in sequence, until it gets an Probe message from B indicating that (a) B has received one of A's messages and, obviously, (b) that B's Probe message gets back to A. B uses the same algorithm, but starts the process from the reception of the first Probe message from A.

Upon changing to a new address pair, transport layer protocol needs to be informed so that it can perform a slow start, or some other form of adaptation to the possibly changed conditions. However, this functionality is outside the scope of REAP and is rather seen as a general multihoming issue.

Similarly, one can also envision that applications would be able to tell the IP or transport layer that the current connection is unsatisfactory and an exploration for a better one would be desirable. This would require an API to be developed, however. In

any case, this is another issue that we treat as being outside the scope of pure address exploration.

5.3. Exploration Order

The exploration process assumes an ability to pick current and alternative address pairs. This process may result in a combinatorial explosion when there are many addresses on both sides, but a back-off procedure is employed to avoid a "signaling storm".

Nodes **MUST** first consult RFC 3484 [6] Section 4 rules to determine what combinations of addresses are allowed from a local point of view, as this reduces the search space. RFC 3484 also provides a priority ordering among different address pairs, making the search possibly faster. Nodes **MAY** also use local information, such as known quality of service parameters or interface types to determine what addresses are preferred over others, and try pairs containing such addresses first. The multihoming protocol also carries preference information in its messages [18].

Discussion note: The preferences may either be learned dynamically or be configured. It is believed, however, that dynamic learning based purely on the multihoming protocol is too hard and not the task this layer should do. Solutions where multiple protocols share their information in a common pool of locators could provide this information from transport protocols, however.

One suggested good implementation strategy is to record the reachability test result (an on/off value) and multiply this by the age of the information. This allows recently tested address pairs to be chosen before old ones.

Out of the set of possible candidate address pairs, nodes **SHOULD** attempt a test through all of them until a working pair is found, and retrying the process as is necessary. However, all nodes **MUST** perform this process sequentially and with exponential back-off. This sequential process is necessary in order to avoid a "signaling storm" when an outage occurs (particularly for a complete site). However, it also limits the number of addresses that can in practice be used for multihoming, considering that transport and application layer protocols will fail if the switch to a new address pair takes too long.

5.4. Protocol Design

REAP is designed as a modular part of SHIM6 in the hopes that it may also be useful in other contexts. This document defines how it is

carried within SHIM6, but the actual protocol messages are self-contained so that it could be carried by other protocols as well.

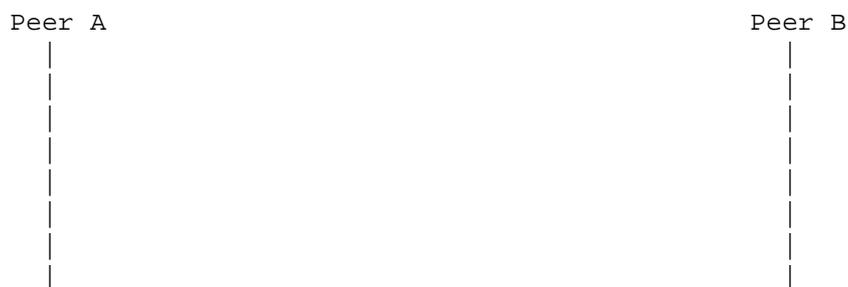
The REAP design allows performing both failure detection and address pair exploration in the same sequence of messages, without a need to designate a specific point when the current address pair is declared inoperational and the search for a new pair begins. This is useful, as the loss of a small number of packets is not a proof that a problem exists. Integrated failure detection and exploration allows us to test multiple address pairs simultaneously, including the current pair in case it starts working again. For instance, the exploration process can refer to keepalive message that succeeded in getting to the peer during the reachability testing phase.

REAP also integrates a return routability function, making it unnecessary to perform another roundtrip before a newly discovered address can be taken into use.

This document defines a minimal set of parameters that are carried by the messages of the protocol. Specifically, we have limited the parameters to those that are necessary to find a working path. We note there may be extensions that are needed in the future for various reasons, such as the desire to support load balancing or finding best paths. An option format has been specified to allow this.

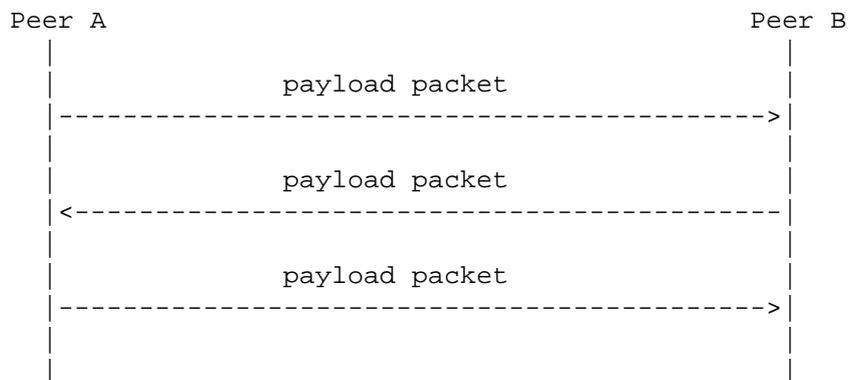
5.5. Example Protocol Runs

This section has examples of REAP protocol runs in typical scenarios. We start with the simplest scenario of two hosts, A and B, that have a SHIM6 connection with each other but are not currently sending any data. As neither side sends anything, they also do not expect anything back, so there are no messages at all:

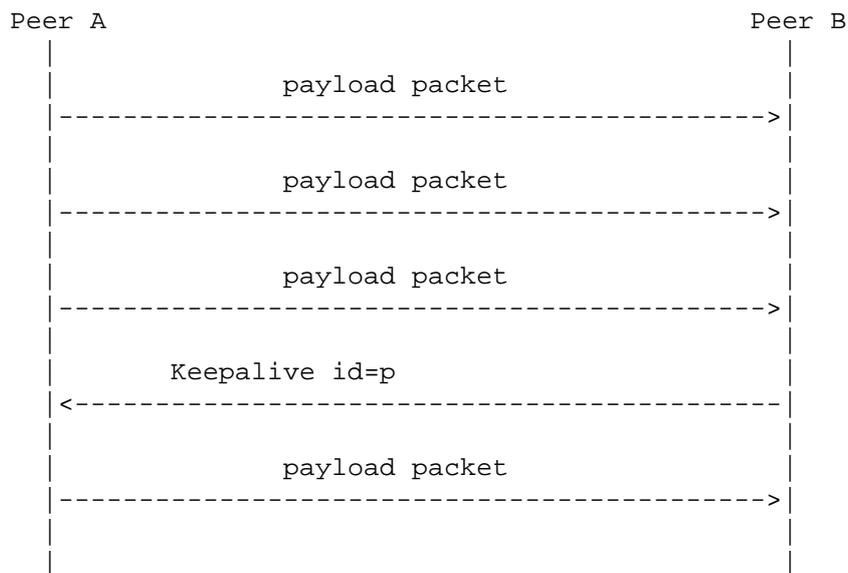


Our second example involves an active connection with bidirectional payload packet flows. Here the reception of data from the peer is taken as an indication of reachability, so again there are no extra

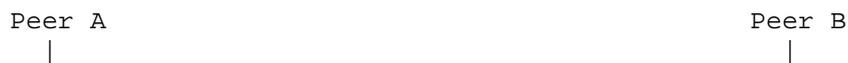
packes:

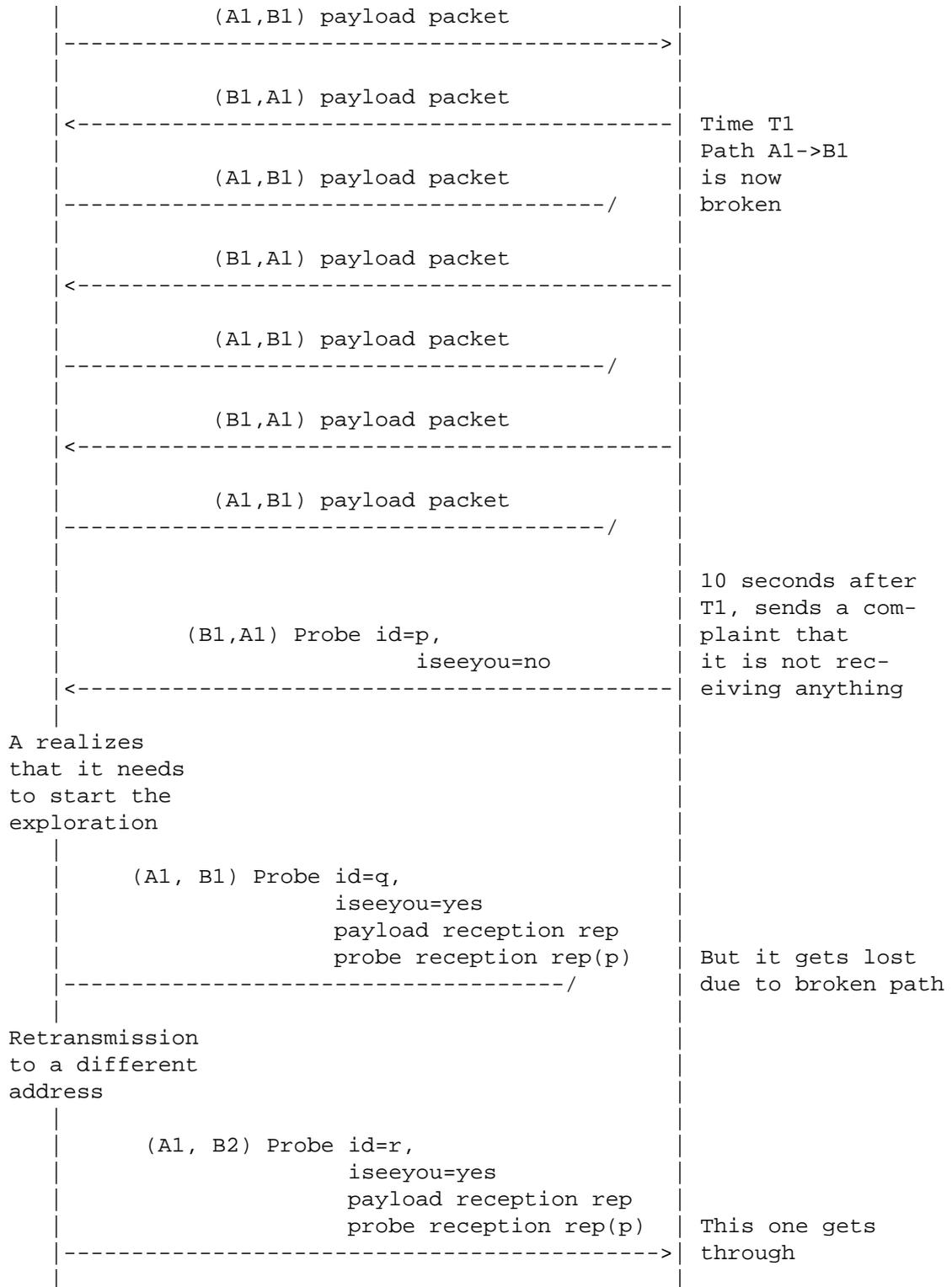


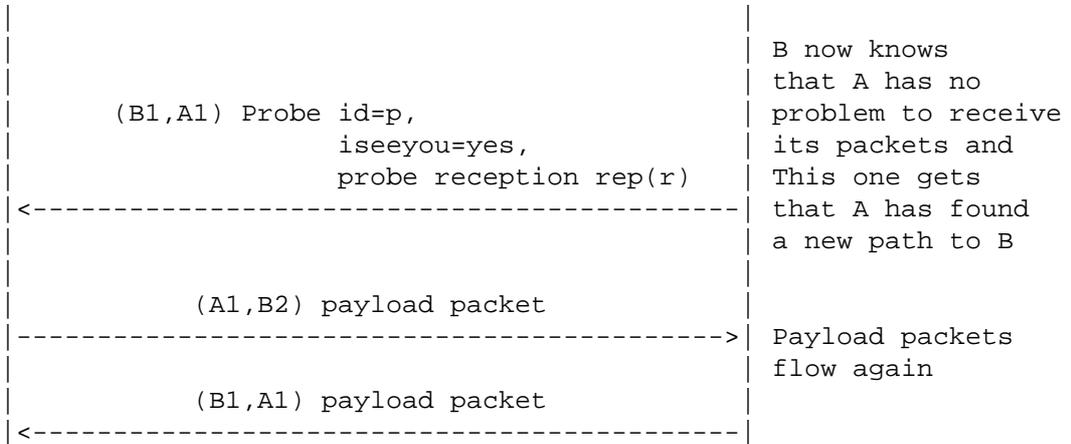
The third example is the first one that involves an actual REAP message. Here the hosts communicate in just one direction, so REAP messages are needed to indicate to the peer that sends payload packets that its packets are getting through:



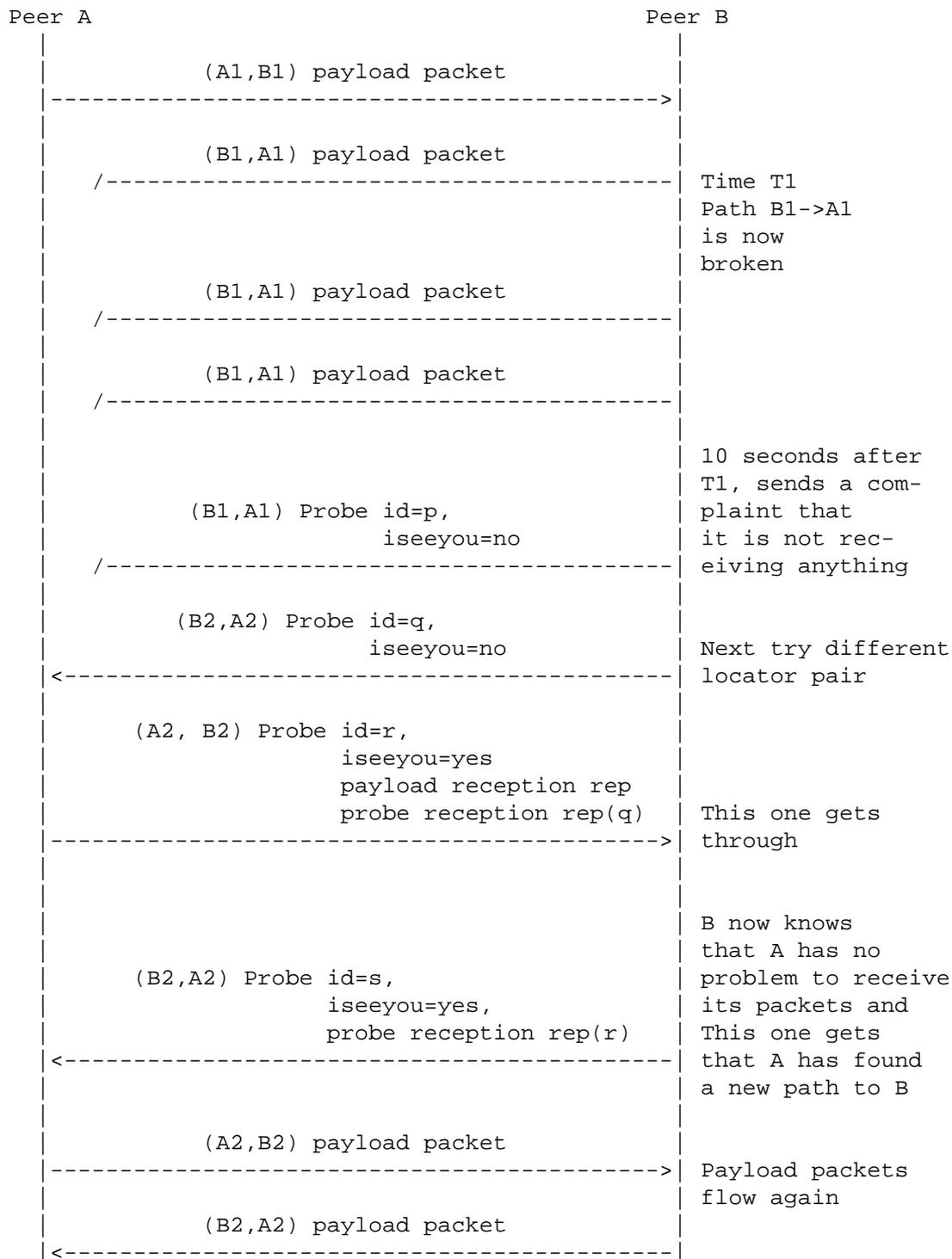
Finally, our last example involves a failure scenario. Here A has addresses A1 and A2 and B has addresses B1 and B2. The currently used address pairs are (A1, B1) and (B1, A1). The first of these becomes broken, which leads to an exploration process:



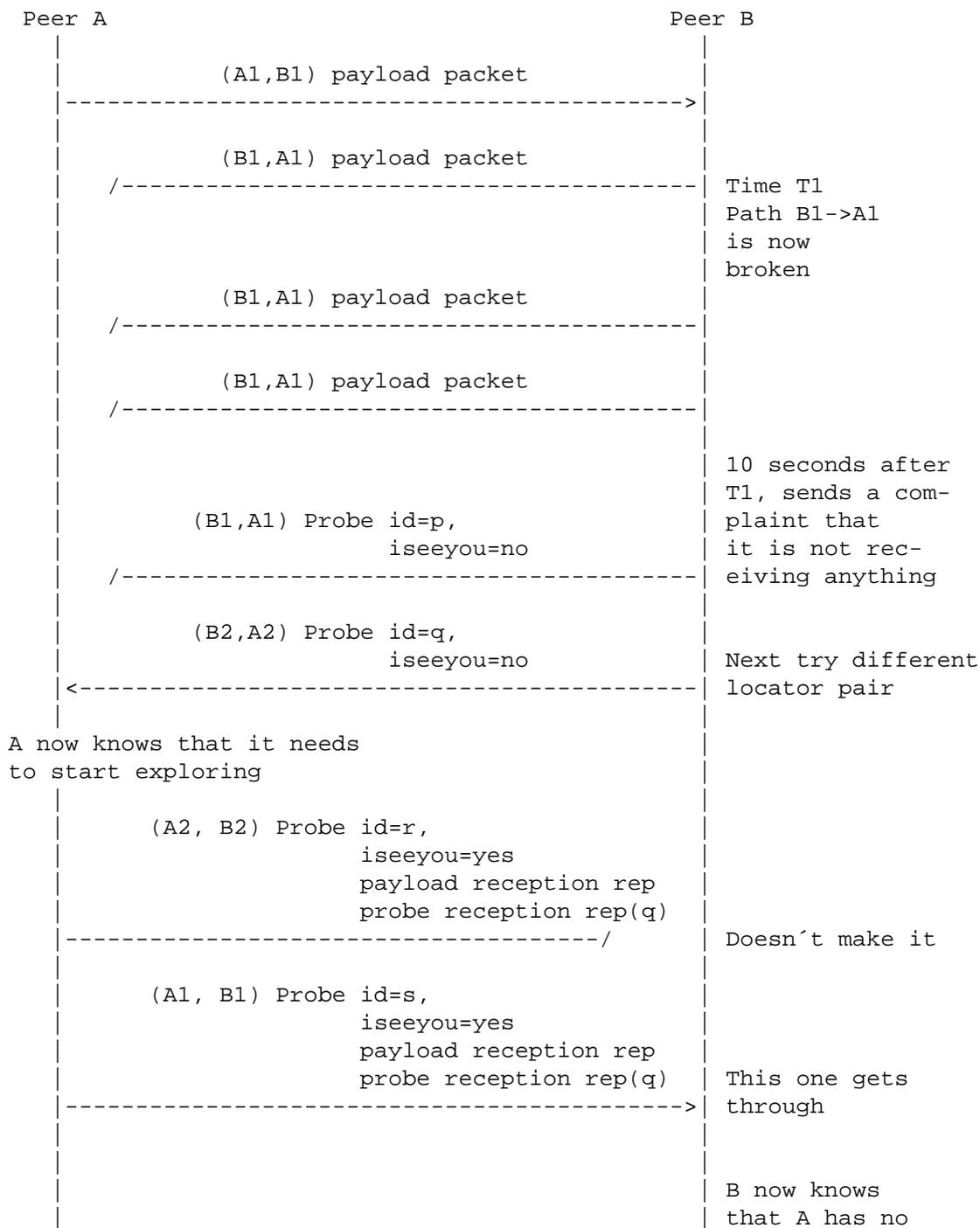


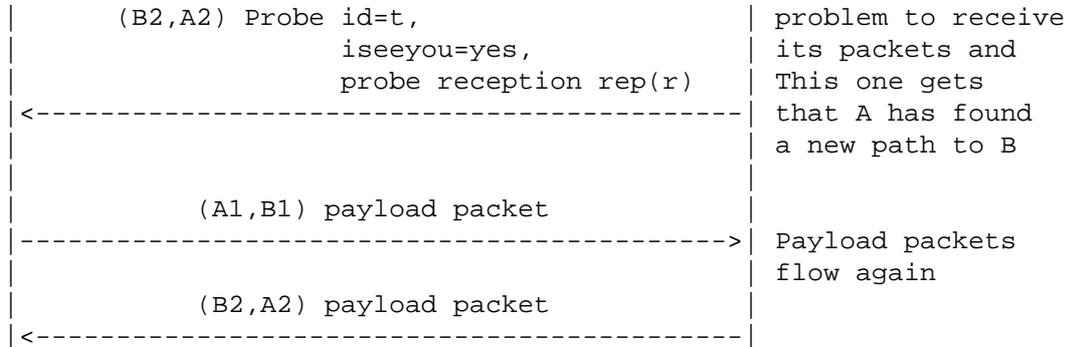


The next example shows when the failure for the current locator pair is in the other direction:



In the next case we have even less luck. The response to the REAP probe doesn't make it in the reverse direction, so both ends end up exploring independently:





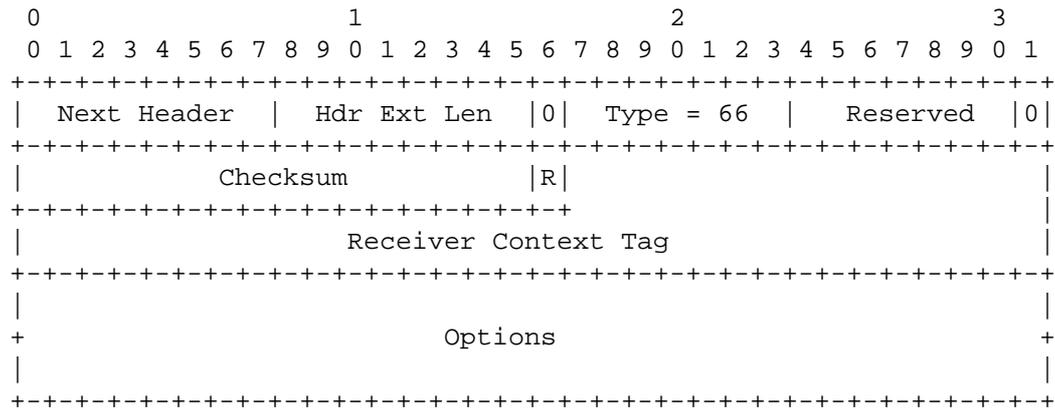
5.6. Limitations

REAP is designed to support failure recovery even in the case of having only unidirectionally operational address pairs. However, due to security concerns discussed in Section 7, the exploration process can typically be run only for a session that has already been established. Specifically, while REAP would in theory be capable of exploration even during connection establishment, its use within the SHIM6 protocol does not allow this.

6. Protocol Definition

6.1. Keepalive Message

The format of the keepalive message is as follows:



Next Header

This value **MUST** be set to NO_NXT_HDR (59).

Type

This field identifies the Probe message and **MUST** be set to 66 (Keepalive).

Reserved

This is a 7-bit field reserved for future use. It is set to zero on transmit, and **MUST** be ignored on receipt.

R

This is a 1-bit field reserved for future use. It is set to zero on transmit, and **MUST** be ignored on receipt.

Receiver Context Tag

This is a 47-bit field for the Context Tag the receiver has allocated for the context.

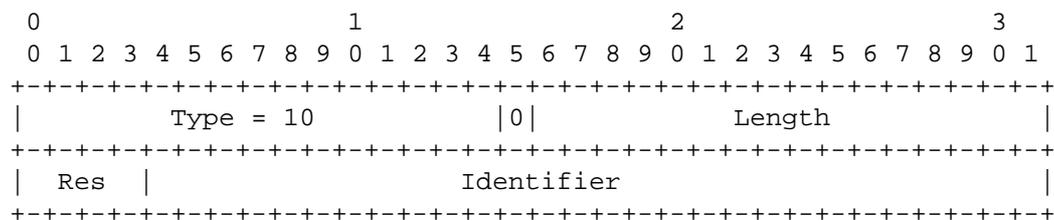
Options

This MUST contain at least the Keepalive option and MAY contain one or more Reachability options. The inclusion of the latter options is not necessary, however, as there are currently no defined options that are useful in a Keepalive message. These options are provided only for future extensibility reasons.

A valid message conforms to the format above, has a Receiver Context Tag that matches to context known by the receiver, is valid shim control message as defined in Section 12.2 of [18], and its shim context state is ESTABLISHED. The receiver processes a valid message by inspecting its options, and executing any actions specified for such options.

The processing rules for this message are the given in more detail in Section 6.4.

6.1.1. Keepalive Option



Type

This value MUST be set to 10 (Keepalive Option).

0

This value MUST be set to 0, as in other SHIM6 options.

Length

This is the length of the option and MUST be calculated as specified in Section 5.14 of [18].

Res

This 4-bit reserved field MUST be set to zero when sending, and ignored on receipt.

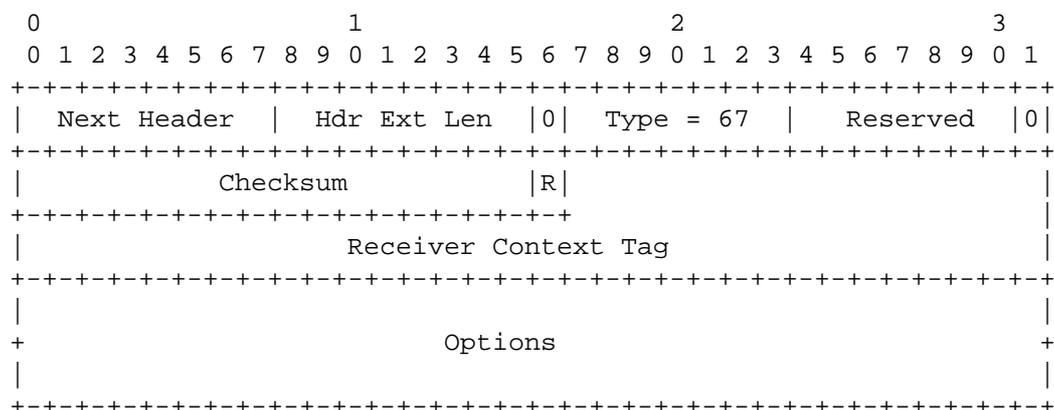
Identifier

This 28-bit field identifies this particular instance of an Keepalive message. This value SHOULD be generated using a random number generator that is known to have good randomness properties [1]. Upon reception, Identifier values from both Keepalive and Probe messages may be copied onto Probe Reception Report options. This allows them to be used for both identifying which packets were received as well as for performing a return routability test.

The processing rules for this option are the given in more detail in Section 6.4.

6.2. Probe Message

This message performs REAP exploration. Its format is as follows:



Next Header

This value MUST be set to NO_NXT_HDR (59).

Type

This field identifies the Probe message and MUST be set to 67 (Probe).

Reserved

This is a 7-bit field reserved for future use. It is set to zero on transmit, and MUST be ignored on receipt.

R

This is a 1-bit field reserved for future use. It is set to zero on transmit, and MUST be ignored on receipt.

Receiver Context Tag

This is a 47-bit field for the Context Tag the receiver has allocated for the context.

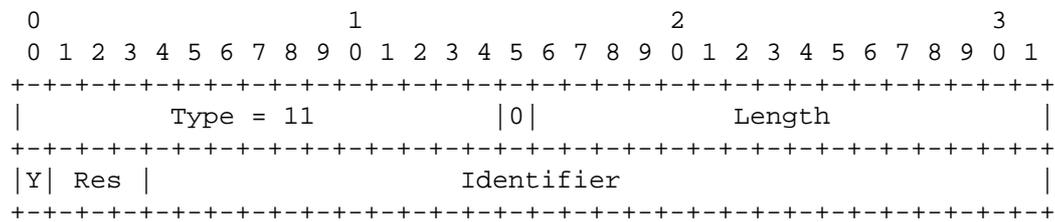
Options

This MUST contain at least the Probe option and MAY contain one or more Reachability options.

A valid message conforms to the format above, has a Receiver Context Tag that matches to a context known by the receiver, is valid shim control message as defined in Section 12.2 of [18], and its shim context state is ESTABLISHED. The receiver processes a valid message by inspecting its options, and executing any actions specified such options. This includes the SHIM6 Probe option found within the options.

The processing rules for this message are the given in more detail in Section 6.4.

6.2.1. Probe Option



Type

This value MUST be set to 11 (Probe Option).

0

This value MUST be set to 0, as in other SHIM6 options.

Length

This is the length of the option and MUST be calculated as specified in Section 5.14 of [18].

Y (The "I See You" flag)

This flag is set to 1 if the sender receives either payload packets or REAP messages from the peer, and 0 otherwise. The determination of when the sender receives something is made during the last Send Timeout seconds (see Section 6.5) when traffic was expected, i.e., when there was either payload traffic or REAP messages.

Upon reception, a value of 1 indicates that the receiver does not need to change its behaviour as the sender is already seeing its packets. A value of 0 indicates that the receiver MUST explore different outgoing address pairs.

Res

This 3-bit reserved field MUST be set to zero when sending, and ignored on receipt.

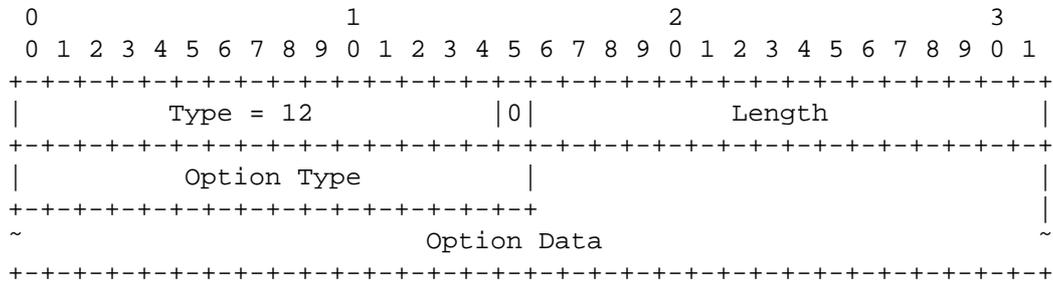
Identifier

This 28-bit field identifies this particular instance of an Probe message. This value SHOULD be generated using a random number generator that is known to have good randomness properties [1]. Upon reception, Identifier values are copied onto Probe Reception Report options. This allows them to be used for both identifying which Probes were received as well as for performing a return routability test.

The processing rules for this option are the given in more detail in Section 6.4.

6.3. Reachability Option

Additional information can be included in Keepalive and Probe messages by the inclusion of the Reachability Options. Their format is as follows:



Type

This value MUST be set to 12 (Reachability option).

0

This value MUST be set to 0, as in other SHIM6 options.

Length

This is the length of the option and MUST be calculated as specified in Section 5.14 of [18].

Option Type

This value identifies the option.

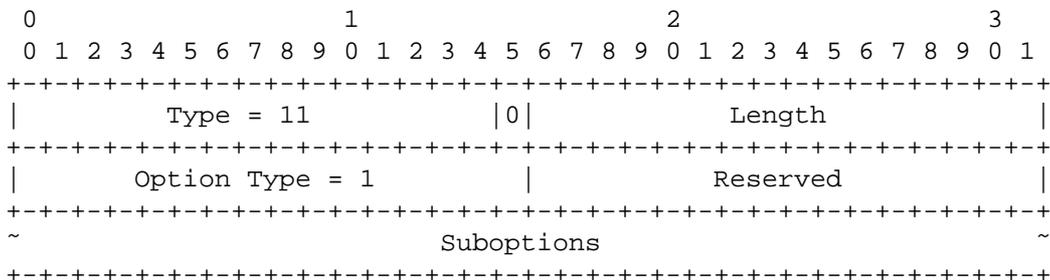
Option Data

Option-specific content.

Unrecognized options MUST be ignored upon receipt. All implementations MUST support the options defined in this specification, however.

6.3.1. Payload Reception Report

This option SHOULD be included in all Probe messages when the sender has recently (within the last Send Timeout seconds) received payload packets from the peer. Its format is as follows:



Type, 0, and Length

These are as specified above.

Reserved

This is a 16-bit field reserved for future use. It is set to zero on transmit, and MUST be ignored on receipt.

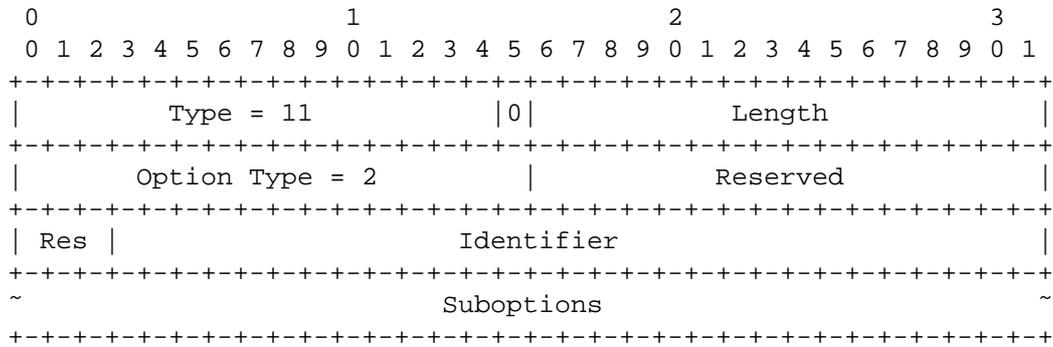
Suboptions

This field is reserved for possible future Reachability options that are carried (recursively) within this option. Unrecognized options MUST be ignored upon receipt. Currently there are no defined options that can be carried here.

6.3.2. Probe Reception Report

This option MUST be included in any Probe message when the sender has recently (within the last Send Timeout seconds) received Probe or Keepalive messages from the peer. Depending on MTU and timing considerations, the sender MAY, however, include options for only some of the received Probe messages. All implementations MUST support sending of at least five such options, however.

The format of this option is as follows:



Type, 0, and Length

These are as specified above.

Option Type

This value identifies the option and MUST be set to 2 (Probe Reception Report).

Reserved

This is a 16-bit field reserved for future use. It is set to zero on transmit, and MUST be ignored on receipt.

Res

This is a 3-bit field reserved for future use. It is set to zero on transmit, and MUST be ignored on receipt.

Identifier

This 32 bit field carries the identifier of the Probe message that was recently received.

Suboptions

This field is reserved for possible future Reachability options that are carried (recursively) within this option. Unrecognized options MUST be ignored upon receipt. Currently there are no defined options that can be carried here.

6.4. Behaviour

The required behaviour of REAP nodes is specified below in the form of a state machine. The externally observable behaviour of an implementation **MUST** conform to this state machine, but there is no requirement that the implementation actually employs a state machine.

On a given context with a given peer, the node can be in one of three states: Operational, Exploring, or ExploringOK. In the Operational state the underlying address pairs are assumed to be operational. In the Exploring state this node has observed a problem and has currently not seen any traffic from the peer. Finally, in the ExploringOK state this node sees traffic from the peer, but peer may not yet see any traffic from this node so that the exploration process needs to continue.

The node maintains also the Send and Keepalive timers. The Send timer reflects the requirement that when this node sends a payload packet there should be some return traffic (either payload packets or Keepalive messages) within Keepalive Timeout seconds. The Keepalive timer reflects the requirement that when this node receives a payload packet there should a similar response towards the peer. The Keepalive timer is only used within the Operational state, and the Send timer in the Operational and ExploringOK states. No timer is running in the Exploring state.

Upon the reception of a payload packet in the Operational state, the node starts the Keepalive timer if it is not yet running, and stops the Send timer if it was running. If the node is in the Exploring state it transitions to the ExploringOK state, sends a Probe message with the I See You flag set to 1 (Yes), and starts the Send timer. In the ExploringOK state the node stops the Send timer if it was running, but does not do anything else. The reception of SHIM6 control messages other than the Keepalive and Probe messages are treated similarly with payload packets.

Upon sending a payload packet in the Operational state, the node stops the Keepalive timer if it was running and starts the Send timer if it was not running. In the Exploring state there is no effect, and in the ExploringOK state the node simply starts the Send timer if it was not yet running. (The sending of SHIM6 control messages is again treated similarly here.)

Upon a timeout on the Keepalive timer the node sends a Keepalive message. This can only happen in the Operational state.

Upon a timeout on the Send timer, the node enters the Exploring state and sends a Probe with I See You set to 0 (No) and stops the

Keepalive timer if it was running.

While in the Exploring state the node keeps retransmitting its Probe messages to different (or same) addresses as defined in Section 5.3. A similar process is employed in the ExploringOk state, except that upon such retransmission the Send timer is started if it was not running already.

Upon the reception of a Keepalive message in the Operational state, the node stops the Send timer, if it was running. If the node is in the Exploring state it transitions to the ExploringOK state, sends a Probe message with the I See You flag set to 1 (Yes), and starts the Send timer. In the ExploringOK state the Send timer is stopped, if it was running.

Upon receiving a Probe with I See You set to 0 (No) the node enters the ExploringOK state, sends a Probe with I See You set to 1 (Yes), stops the Keepalive timer if it was running, and restarts the Send timer.

The behavior upon the reception of a Probe message with I see You set to 1 (Yes) depends on whether it contains a Probe Reception Report that refers to a Probe that this node has sent to the peer such that the I See You was set to 1 in that message. If not, the node sends a Probe message with I See You set to 1 (Yes), restarts the Send timer, stops the Keepalive timer if it was running, and transitions to the Operational state.

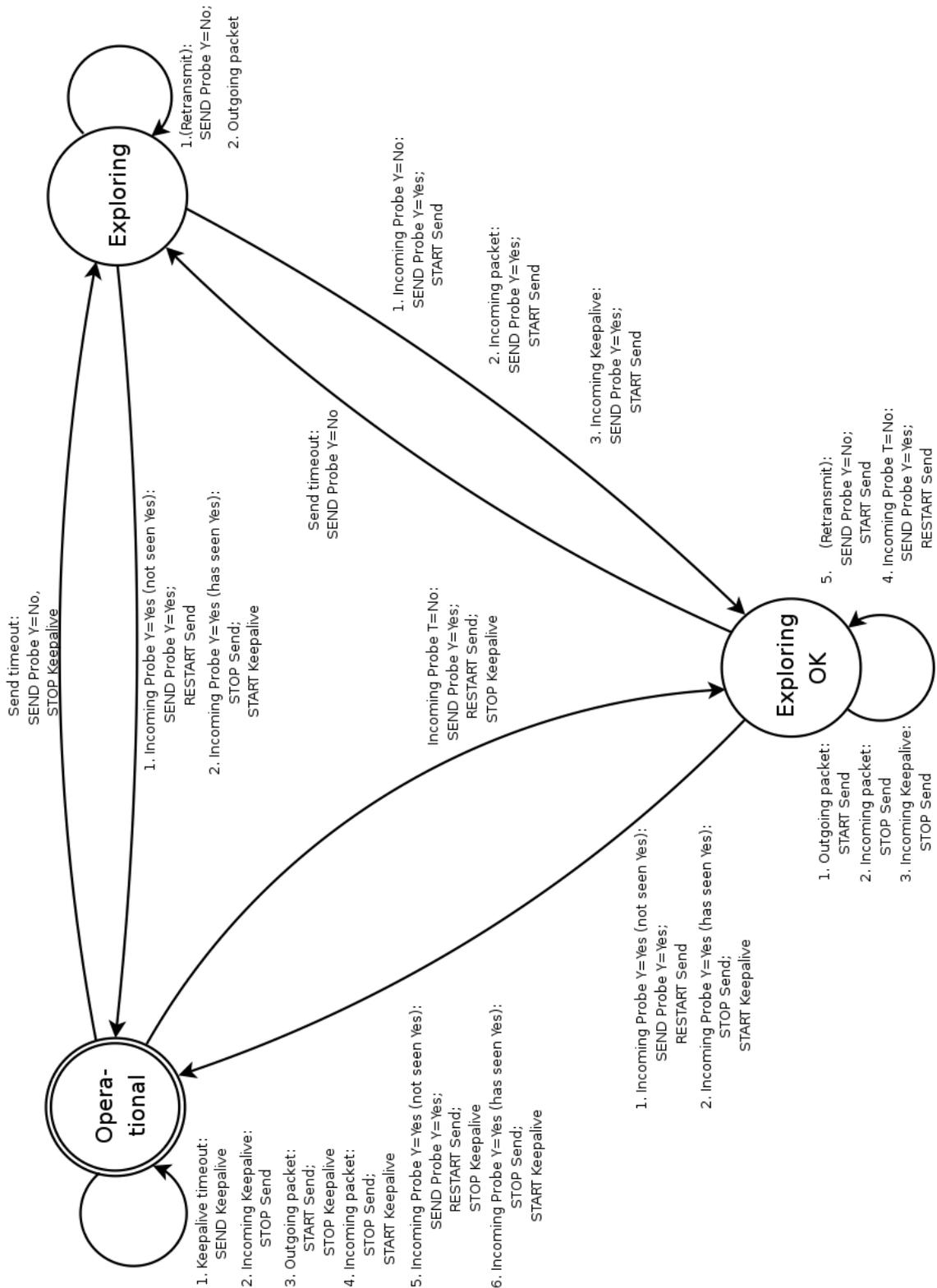
If there was no such Probe Reception Report, the stops the Send timer if it was running, starts the Keepalive timer if it was not yet running, and transitions to the Operational state.

Note: This check is necessary in order to terminate the exploration process when both parties are happy and know that their peers are happy as well.

The reachability detection and exploration process has no effect on payload communications until a new working address pairs have actually been confirmed. Prior to that the payload packets continue to be sent to the previously used addresses.

Garbage collection of SHIM6 contexts terminates contexts that are either unused or have failed due to the inability of the exploration process to find a working pair.

In the PDF version of this specification, an informational drawing illustrates the state machine. Where the text and the drawing differ, the text takes precedence.



A tabular representation of the state machine is shown below. Like the drawing, this representation is only informational.

1. EVENT: Incoming payload packet

=====

Operational	Exploring	ExploringOk
STOP Send; START Keepalive	SEND Probe Y=Yes; START Send; GOTO ExploringOk	STOP Send

2. EVENT: Outgoing payload packet

=====

Operational	Exploring	ExploringOk
START Send; STOP Keepalive	-	START Send

3. EVENT: Keepalive timeout

Operational	Exploring	ExploringOk
SEND Keepalive	-	-

4. EVENT: Send timeout

=====

Operational	Exploring	ExploringOk
SEND Probe Y=No; STOP Keepalive; GOTO EXPLORING	-	SEND Probe Y=No GOTO EXPLORING

5. EVENT: Reception of the Keepalive message

=====

Operational	Exploring	ExploringOk
STOP Send	SEND Probe Y=Yes;	STOP Send

START Send;
GOTO ExploringOk

6. EVENT: Reception of the Probe message with Y=No
=====

Operational	Exploring	ExploringOk
-----	-----	-----
SEND Probe Y=Yes	SEND Probe Y=Yes;	SEND Probe Y=Yes;
STOP Keepalive;	START Send;	RESTART Send
RESTART Send;	GOTO EXPLORINGOK	
GOTO EXPLORINGOK		

7. EVENT: Reception of the Probe message with Y=Yes
(peer reports not seeing yet a Probe with Y=Yes)
=====

Operational	Exploring	ExploringOk
-----	-----	-----
SEND Probe Y=Yes;	SEND Probe Y=Yes;	SEND Probe Y=Yes;
RESTART Send;	RESTART Send;	RESTART Send;
STOP Keepalive	GOTO OPERATIONAL	GOTO OPERATIONAL

8. EVENT: Reception of the Probe message with Y=Yes
(peer reports seeing a Probe with Y=Yes)
=====

Operational	Exploring	ExploringOk
-----	-----	-----
STOP Send	STOP Send;	STOP Send;
START Keepalive	START Keepalive	START Keepalive
	GOTO OPERATIONAL	GOTO OPERATIONAL

9. EVENT: Retransmission
=====

Operational	Exploring	ExploringOk
-----	-----	-----
-	SEND Probe Y=No	SEND Probe Y=Yes
		START Send

6.5. Protocol Constants

The following protocol constants are defined:

Send Timeout	10 seconds
Keepalive Timeout	3 seconds

7. Security Considerations

Attackers may spoof various indications from lower layers and the network in an effort to confuse the peers about which addresses are or are not working. For example, attackers may spoof ICMP error messages in an effort to cause the parties to move their traffic elsewhere or even to disconnect. Attackers may also spoof information related to network attachments, router discovery, and address assignments in an effort to make the parties believe they have Internet connectivity when in reality they do not.

This may cause use of non-preferred addresses or even denial-of- service.

This protocol does not provide any protection of its own for indications from other parts of the protocol stack. However, this protocol has weak resistance against incorrect information from these sources in the sense that it performs its own tests prior to picking a new address pair. Denial-of- service vulnerabilities remain, however, as do vulnerabilities against on path attackers.

Some aspects of these vulnerabilities can be mitigated through the use of techniques specific to the other parts of the stack, such as properly dealing with ICMP errors [22], link layer security, or the use of [12] to protect IPv6 Router and Neighbor Discovery.

This protocol is designed to be used in situations where other parts of the stack have ensured that a set of addresses belong together, such as via SHIM6 HBAs [17]. That is, REAP itself provides no assurance that a set of addresses belongs to the same host. Similarly, REAP provides only minimal protection against third party flooding attacks; when REAP is run its Probe identifiers can be used as a return routability check that the claimed address is indeed willing to receive traffic. However, this needs to be complemented with another mechanism to ensure that the claimed address is also the correct host. In SHIM6 this is performed by binding all operations to context tags.

Finally, the exploration itself can cause a number of packets to be sent. As a result it may be used as a tool for packet amplification in flooding attacks. In order to prevent this it is required that the protocol employing REAP has built-in mechanisms to prevent this. For instance, in SHIM6 contexts are created only after a relatively large number of packets has been exchanged, a cost which reduces the attractiveness of using SHIM6 and REAP for amplification attacks. However, such protections are typically not present at connection establishment time. When exploration would be needed for connection establishment to succeed, its usage would result in an amplification

vulnerability. As a result, SHIM6 does not support the use of REAP in connection establishment stage.

8. IANA Considerations

This document creates one new name spaces under the new SHIM6 Reachability Protocol repository. The name space is for Reachability Option Type (Section 6.3) and it has one reserved value (0) and two defined values, 1 (Payload Reception Report defined in Section 6.3.1) and 2 (Probe Reception Report defined in Section 6.3.2). Further allocations within this 16-bit field can be made through Specification Required. The range from 65000 to 65535 is reserved for experimental use.

9. References

9.1. Normative References

- [1] Eastlake, D., Crocker, S., and J. Schiller, "Randomness Recommendations for Security", RFC 1750, December 1994.
- [2] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [3] Narten, T., Nordmark, E., and W. Simpson, "Neighbor Discovery for IP Version 6 (IPv6)", RFC 2461, December 1998.
- [4] Thomson, S. and T. Narten, "IPv6 Stateless Address Autoconfiguration", RFC 2462, December 1998.
- [5] Droms, R., Bound, J., Volz, B., Lemon, T., Perkins, C., and M. Carney, "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)", RFC 3315, July 2003.
- [6] Draves, R., "Default Address Selection for Internet Protocol version 6 (IPv6)", RFC 3484, February 2003.
- [7] Choi, J., "Detecting Network Attachment in IPv6 Goals", draft-ietf-dna-goals-00 (work in progress), June 2004.
- [8] Moore, N., "Optimistic Duplicate Address Detection for IPv6", draft-ietf-ipv6-optimistic-dad-01 (work in progress), June 2004.
- [9] Hinden, R. and B. Haberman, "Unique Local IPv6 Unicast Addresses", draft-ietf-ipv6-unique-local-addr-05 (work in progress), June 2004.

9.2. Informative References

- [10] Stewart, R., Xie, Q., Morneault, K., Sharp, C., Schwarzbauer, H., Taylor, T., Rytina, I., Kalla, M., Zhang, L., and V. Paxson, "Stream Control Transmission Protocol", RFC 2960, October 2000.
- [11] Rosenberg, J., Weinberger, J., Huitema, C., and R. Mahy, "STUN - Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs)", RFC 3489, March 2003.
- [12] Arkko, J., Kempf, J., Zill, B., and P. Nikander, "SEcure Neighbor Discovery (SEND)", RFC 3971, March 2005.
- [13] Aboba, B., "Detection of Network Attachment (DNA) in IPv4",

- draft-ietf-dhc-dna-ipv4-08 (work in progress), July 2004.
- [14] Nikander, P., "End-Host Mobility and Multi-Homing with Host Identity Protocol", draft-ietf-hip-mm-00 (work in progress), October 2004.
- [15] Eronen, P., "IKEv2 Mobility and Multihoming Protocol (MOBIKE)", draft-ietf-mobike-protocol-03 (work in progress), September 2005.
- [16] Rosenberg, J., "Interactive Connectivity Establishment (ICE): A Methodology for Network Address Translator (NAT) Traversal for Multimedia Session Establishment Protocols", draft-ietf-mmusic-ice-02 (work in progress), July 2004.
- [17] Bagnulo, M., "Hash Based Addresses (HBA)", draft-ietf-shim6-hba-00 (work in progress), July 2005.
- [18] Nordmark, E., "Level 3 multihoming shim protocol", draft-ietf-shim6-proto-00 (work in progress), October 2005.
- [19] Beijnum, I., "Shim6 Reachability Detection", draft-ietf-shim6-reach-detect-00 (work in progress), July 2005.
- [20] Stewart, R., "Stream Control Transmission Protocol (SCTP) Dynamic Address Reconfiguration", draft-ietf-tsvwg-addip-sctp-10 (work in progress), January 2005.
- [21] Bagnulo, M., "Address selection in multihomed environments", draft-bagnulo-shim6-addr-selection-00 (work in progress), October 2005.
- [22] Gont, F., "ICMP attacks against TCP", draft-gont-tcpm-icmp-attacks-00 (work in progress), August 2004.
- [23] Huitema, C., "Address selection in multihomed environments", draft-huitema-multi6-addr-selection-00 (work in progress), October 2004.
- [24] Rosenberg, J., "Traversal Using Relay NAT (TURN)", draft-rosenberg-midcom-turn-05 (work in progress), July 2004.
- [25] Aura, T., Roe, M., and J. Arkko, "Security of Internet Location Management", In Proceedings of the 18th Annual Computer Security Applications Conference, Las Vegas, Nevada, USA., December 2002.

Appendix A. Contributors

This draft attempts to summarize the thoughts and unpublished contributions of many people, including the MULTI6 WG design team members Marcelo Bagnulo Braun, Iljitsch van Beijnum, Erik Nordmark, Geoff Huston, Margaret Wasserman, and Jukka Ylitalo, the MOBIKE WG contributors Pasi Eronen, Tero Kivinen, Francis Dupont, Spencer Dawkins, and James Kempf, and my colleague Pekka Nikander at Ericsson. This draft is also in debt to work done in the context of SCTP [10] and HIP [14].

Appendix B. Acknowledgements

The author would also like to thank Christian Huitema, Pekka Savola, and Hannes Tschofenig for interesting discussions in this problem space, and for their comments on earlier versions of this draft.

Authors' Addresses

Jari Arkko
Ericsson
Jorvas 02420
Finland

Email: jari.arkko@ericsson.com

Iljitsch van Beijnum
Muada
The Netherlands

Email: iljitsch@muada.com

Intellectual Property Statement

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

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

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

Disclaimer of Validity

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

Copyright Statement

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

Acknowledgment

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