

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: January 7, 2016

C. Cassar
Cisco Systems
I. Kouvelas
Arista Networks Inc.
D. Lewis
J. Arango
J. Leong
Cisco Systems
July 6, 2015

LISP Map Server Reliable Transport
draft-kouvelas-lisp-map-server-reliable-transport-00.txt

Abstract

The communication between LISP ETRs and Map-Servers is based on unreliable UDP message exchange coupled with periodic message transmission in order to maintain soft state. The drawback of periodic messaging is the constant load imposed on both the ETR and the Map-Server. New use cases for LISP have increased the amount of state that needs to be communicated with requirements that are not satisfied by the current mechanism. This document introduces the use of a reliable transport for ETR to Map-Server communication in order to eliminate the periodic messaging overhead, while providing reliability, flow-control and endpoint liveness detection.

This document has been renamed to avoid ambiguity. It is an update to [I-D.kouvelas-lisp-reliable-transport].

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). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

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."

This Internet-Draft will expire on January 7, 2016.

Copyright Notice

Copyright (c) 2015 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. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1.	Introduction	2
2.	Requirements Notation	3
3.	Message Format	4
4.	Session Establishment	5
5.	Error Notifications	5
6.	EID Prefix Registration	7
6.1.	Reliable Mapping Registration Messages	7
6.1.1.	Registration Message	8
6.1.2.	Registration Acknowledgement Message	8
6.1.3.	Registration Rejected Message	9
6.1.4.	Registration Refresh Message	10
6.1.5.	Mapping Notification Message	10
6.2.	ETR Behavior	11
6.3.	Map-Server Behavior	15
7.	Security Considerations	16
8.	IANA Considerations	16
8.1.	LISP Reliable Transport Message Types	16
8.2.	Transport Protocol Port Numbers	16
9.	Acknowledgments	16
10.	Normative References	17
	Authors' Addresses	17

1. Introduction

The communication channel between LISP ETRs and Map-Servers is based on unreliable UDP message exchange [RFC6833]. Where required, reliability is pursued through periodic retransmissions that maintain soft state on the peer. Map-Register messages are retransmitted every minute by an ETR and the Map-Server times out its state if the state is not refreshed for three successive periods. When registering multiple EID-Prefixes, the ETR includes multiple mapping

records in the Map-Register message. Packet size limitations provide an upper bound to the number of mapping records that can be placed in each Map-Register message. When the ETR has more EID-Prefixes to register than can be packed in a single Map-Register message, the mapping records for the EID-Prefixes are split across multiple Map-Register messages.

The drawback of the periodic registration is the constant load that it introduces on both the ETR and the Map-Server. The ETR uses resources to periodically build and transmit the Map-Register messages, and to process the resulting Map-Notify messages issued by the Map-Server. The Map-Server uses resources to process the received Map-Register messages, update the corresponding registration state, and build and transmit the matching Map-Notify messages. When the number of EID-Prefixes to be registered by an ETR is small, the resulting load imposed by periodic registrations may not be significant. The ETR will only transmit a single Map-Register message each period that contains a small number of mapping records.

In some LISP deployments, a large set of EID-Prefixes must be registered by each ETR (e.g. mobility, database redistribution). Use cases with a large set of EID-Prefixes behind an ETR will result in a much higher load. An example is LISP mobility deployments where EID-Prefixes are limited to host entries. ETRs may have thousands of hosts to register resulting in hundreds of Map-Register and Map-Notify messages per registration period.

A transport is required for the ETR to Map-Server communication that provides reliability, flow-control and endpoint liveness notifications. This document describes the use of TCP or SCTP as a LISP reliable transport. The initial application for the LISP reliable transport session is the support of scalable EID prefix registration. The reliable session mechanism is defined to be extensible so that it can support additional LISP communication requirements as they arise using a single reliable transport session between an ETR and a Map-Server. The use of the reliable transport session for EID prefix registration is an alternative and does not replace the existing UDP based mechanism.

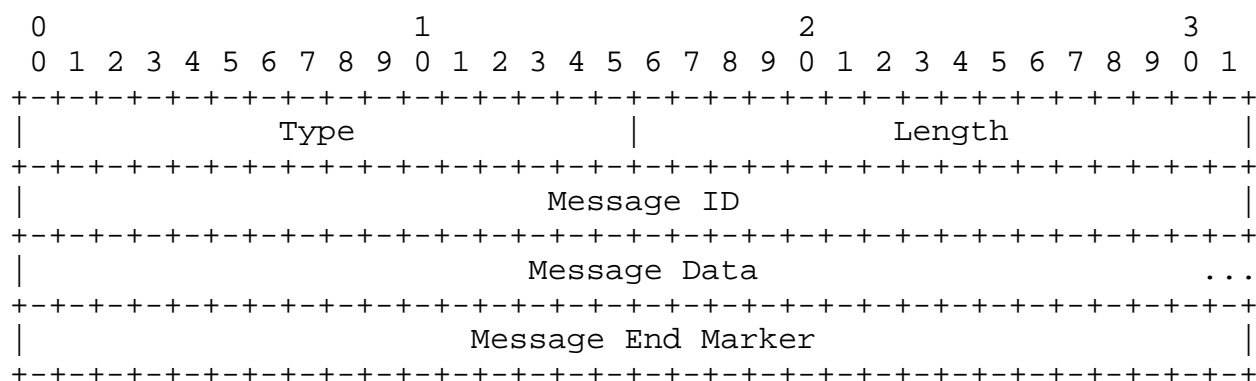
2. Requirements Notation

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].

3. Message Format

A single LISP reliable transport session may carry information for multiple LISP applications. One such application is the registration of EID to RLOC mappings that operates over a session between an ETR and a Map-Server. Communication over a session is based on the exchange of messages. This document defines a base set of messages to support session establishment and management. It also defines the messages for the EID to RLOC mapping registration application.

To support protocol extensibility when new applications, or extensions to existing applications are introduced, the messages are based on a TLV format.



Reliable transport message format

- o Type: 16 bit type field identifying the message type.
- o Length: 16 bit field that provides the total size of the message in octets including the length, type and end marker fields. The length allows the receiver to locate the next message in the TCP stream. The minimum value of the length field is 8.
- o ID: A 32-bit value that identifies the message. May be used by the receiver to identify the message in replies or notification messages.
- o Data: Type specific message contents.
- o End Marker: A 32-bit message end marker that must be set to 0x9FACADE9. The End Marker is used by the receiver to validate that it has correctly parsed or skipped a message and provides a method to detect formatting errors. Note that message data may also contain this marker, and that the marker itself is not sufficient for parsing the message.

The base message format does not indicate how the peer should deal with the message in cases where the message type is not supported/understood. This is best dealt with by the application. For example, in case an error notification is returned, or an expected acknowledgement message is not received, the application might choose various courses of action; from simply logging that the feature is not supported, all the way to tearing the relationship with the peer down for the feature, or for all LISP features.

4. Session Establishment

To ensure backwards compatibility, the map server and ETR MUST communicate via unreliable UDP messages until a TCP session between the two is successfully established.

The map server authenticates the ETR with the authentication data contained in the first UDP map-register message it receives from the ETR. Once the ETR is authenticated, the map server performs a passive open by listening on TCP port 4342, and does not qualify the remote port. As a security measure, the map server accepts TCP connections only from those ETRs that have been authenticated via UDP map-register messages.

The ETR assumes the active role of the TCP session establishment by connecting to the map server once it has received a UDP map-notify message.

When a TCP session goes down, UDP authentication must take place before a new TCP session is established. The map-server will not accept a connection from the ETR until a UDP map-register has been received. Similarly, the ETR will not attempt to establish a session with the map server until an UDP map-notify message has been received.

A single reliable transport session is established between the map server and the ETR to cover all communication needs. For example, an ETR that has EID prefix registrations for multiple EID instances and EID address families will only establish a single session with the map server.

5. Error Notifications

The error notification message is used to communicate base reliable transport session communication errors. LISP applications making use of the reliable transport session and having to communicate application specific errors must define their own messages to do so. An error notification is issued when the receiver of a message does not recognize the message type or cannot parse the message contents.

The notification includes the offending message type and ID and as much of the offending message data as the notification sender wishes to.



Error notification message format

- o Error Code: An 8 bit field identifying the type of error that occurred. Defined errors are:
 - * Unrecognized message type.
 - * Message format error.
- o Reserved: Set to zero by the sender and ignored by the receiver.
- o Offending Message Type: 16 bit type field identifying the message type of the offending message that triggered this error notification. This is copied from the Type field of the offending message.
- o Offending Message Length: 16 bit field that provides the total size of the offending message in octets. This is copied from the Length field of the offending message.
- o Offending Message ID: A 32-bit field that is set to the Message ID field of the offending message.
- o Offending Message Data: The Data from the offending message that triggered this error notification. The sender of the notification may include as much of the original data as is deemed necessary.

The length of the Offending Message Data field is not provided by the Offending Message Length field and is determined by subtracting the size of the other fields in the message from the Length field. It is valid to not include any of the offending message data when sending an error notification.

- o End Marker: A 32-bit message end marker that must be set to 0x9FACADE9. The End Marker is used by the receiver to validate that it has correctly parsed or skipped a message and provides a method to detect formatting errors. Note that message data may also contain this marker, and that the marker itself is not sufficient for parsing the message.

An error notification cannot be the offending message in another error notification and MUST NOT trigger such a message.

6. EID Prefix Registration

EID prefix registration uses the reliable transport session between an ETR and a Map-Server to communicate the ETR local EID database EID to RLOC mappings to the Map-Server. In contrast to the UDP based periodic registration, mapping information over the reliable transport session is only sent when there is new information available for the Map-Server. The Map-Server does not maintain a timer to expire registrations communicated over the reliable transport session. Instead an explicit de-registration (a registration carrying a zero TTL) is needed to delete the state maintained by the Map-Server.

The key used to identify registration mapping records in the ETR to Map-Server communication is the EID prefix. The prefix may be specified using an LCAF encoding that includes an EID instance ID.

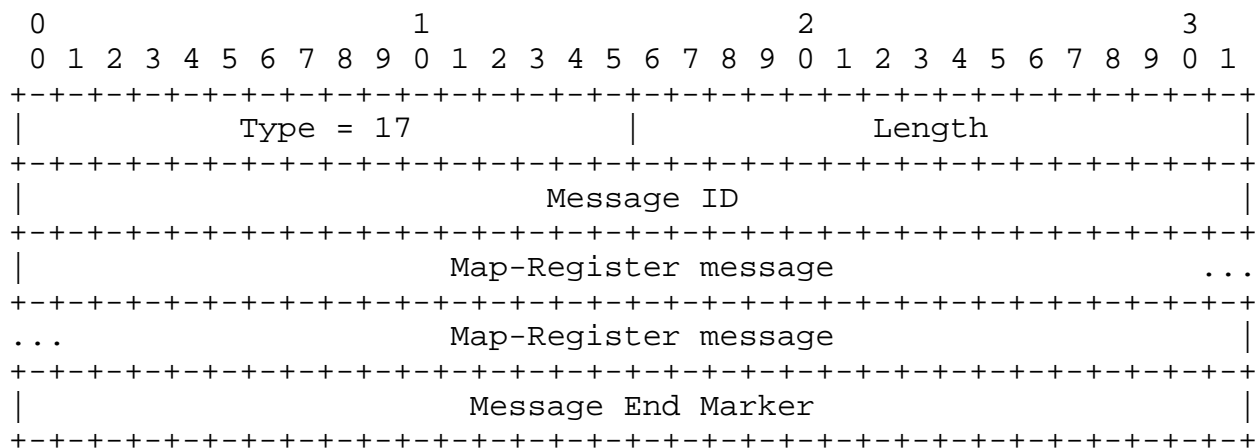
When the reliable transport session goes down, registration mappings learned by the Map-Server are treated as periodic UDP registrations and a timer is used to expire them after 3 minutes. During this period UDP based registrations or the re-establishment of the reliable transport session and subsequent communication of a new mapping can update the EID prefix mapping state.

6.1. Reliable Mapping Registration Messages

This section defines the LISP reliable transport session messages used to communicate local EID database registrations between the ETR and the Map-Server.

6.1.1. Registration Message

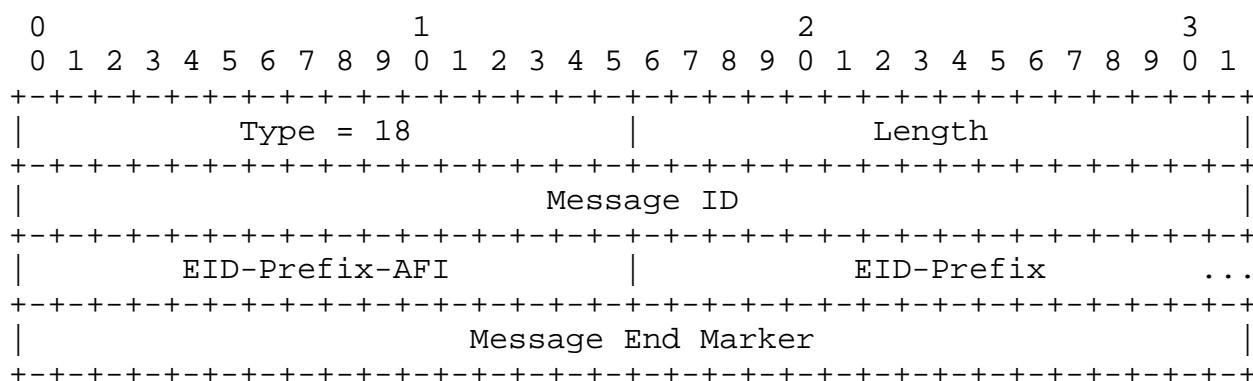
The reliable transport Registration message is used to communicate EID to RLOC mapping registrations from the ETR to the Map-Server. The Registration message uses exactly the same format as the UDP Map-Register message but instead of the IP/UDP header, the Map-Register is placed within the value section of the reliable transport TLV. A common message format is proposed to leverage the authentication features built into the UDP Map-Register message and increase code reuse.



Registration message format

6.1.2. Registration Acknowledgement Message

The Acknowledgement message is sent from the Map-Server to the ETR to confirm successful registration of an EID prefix previously communicated by a reliable transport session Registration message. The Registration Acknowledgement message does not carry a mapping record (the map servers view of the mapping). This is accomplished by the LISP reliable transport Map Notification message.

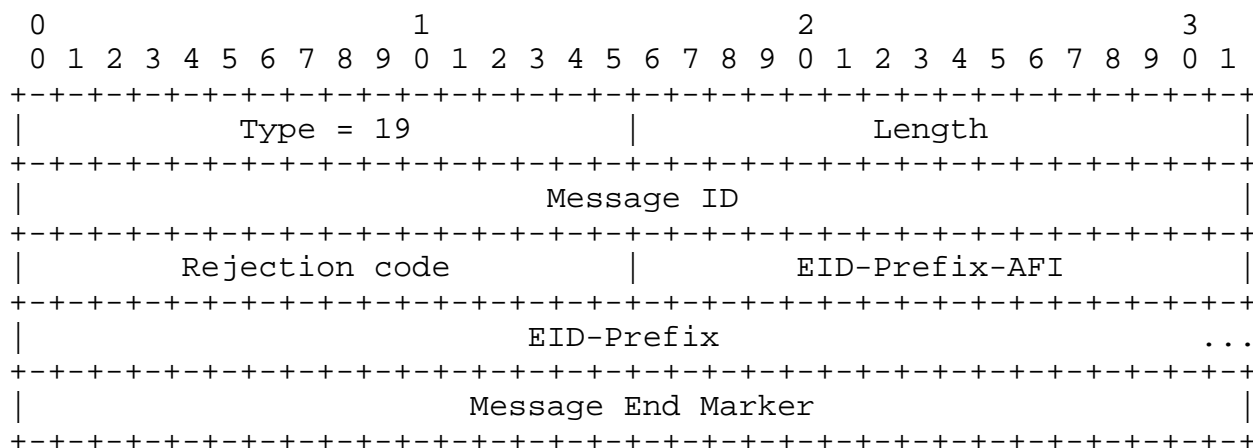


Registration Acknowledgement message format

- o EID-Prefix AFI: Address family identifier for the EID prefix in the following field.
- o EID-Prefix: The EID prefix from the received Registration.

6.1.3. Registration Rejected Message

Negative acknowledgement sent from the Map-Server to the ETR to indicate that the registration of a specific EID prefix was rejected. The ETR must keep track of the fact that the registration of the EID prefix was rejected by the Map-Server and be prepared to re-register the mapping when requested through a failed Registration Refresh request.



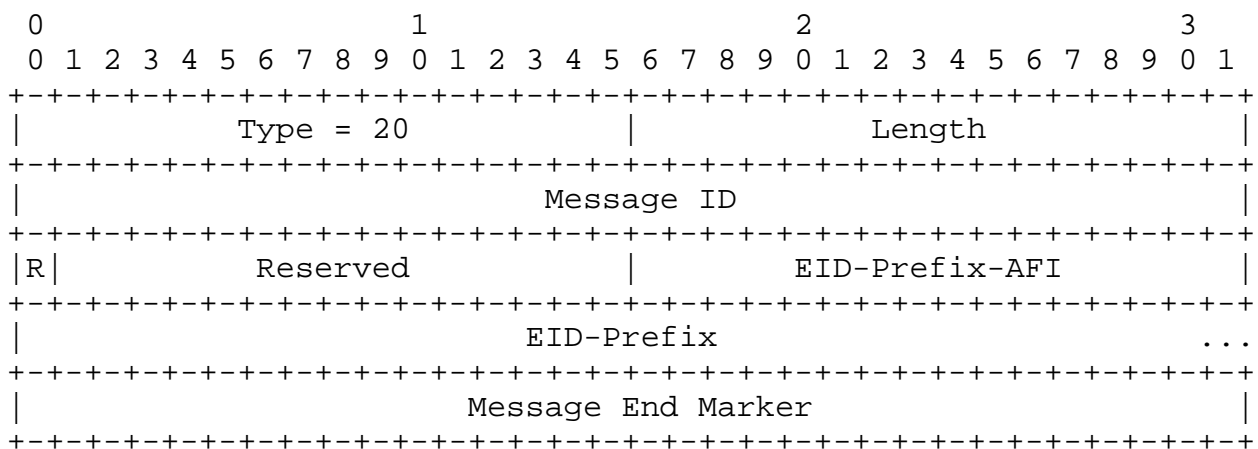
Registration Rejected message format

- o Rejection code: Code identifying the reason for which the Map-Server rejected the registration. Codes:
 - * 1 - Not a valid site EID prefix.

- * 2 - Authentication failure.
- * 3 - Locator set not allowed.
- o EID-Prefix AFI: Address family identifier for the EID prefix in the following field.
- o EID-Prefix: The EID prefix from the received Registration.

6.1.4. Registration Refresh Message

Sent by the Map-Server to the ETR to request the re-transmission of EID prefix database mapping Registration messages.

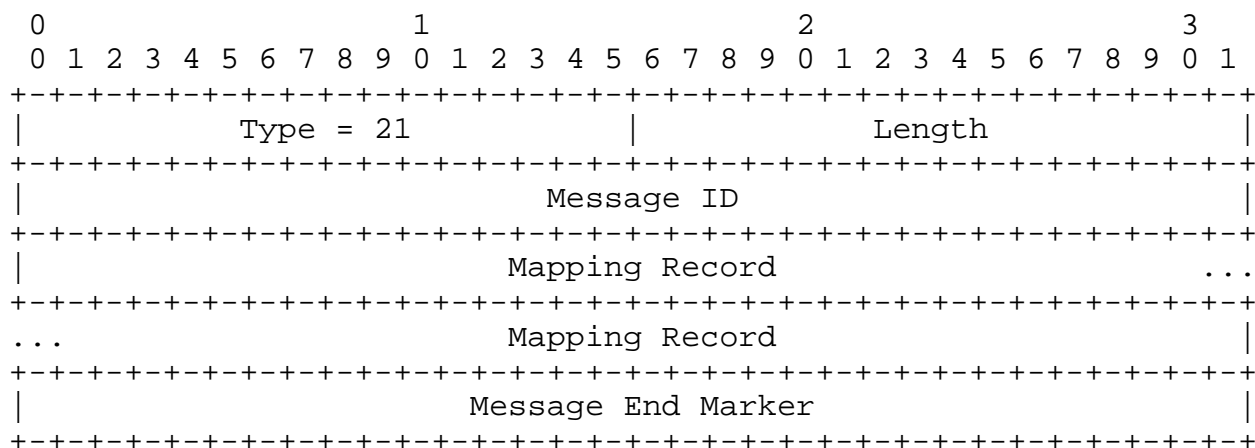


Registration Refresh message format

- o R: Request from the ETR to only refresh registrations that have been previously rejected by the Map-Server.
- o EID prefix, and its more specifics, to refresh. The prefix can be in LCAF format allowing specification of a complete refresh (unspecified prefix), refresh of all the prefixes under an EID instance or even of more specific registrations under a specific EID prefix.

6.1.5. Mapping Notification Message

Mapping Notification messages communicate the Map-Server view of the mapping for an EID prefix and no longer serve as a registration acknowledgement. Mapping Notifications do not need message level authentication as they are received over a reliable transport session to a known Map-Server. Note that reliable transport Mapping Notification messages do not reuse the UDP Map-Notify message format.



Mapping Notification message format

6.2. ETR Behavior

The ETR operates the following per EID prefix, per MS state machine that defines the reliable transport EID prefix registration behavior.

There are five states:

- o No state: The local EID database prefix does not exist.
- o Periodic: The local EID database prefix is being periodically registered through UDP Map-Register messages as specified in [1].
- o Stable: From the ETR's perspective, no registrations are due to be sent to the peer. The session to the peer is up, and the peer has either acknowledged the registration, or is expected to request a refresh in the future.
- o AckWait: A Registration message for the prefix has been transmitted to the Map-Server and the ETR is waiting for either a Registration Acknowledge or Registration Rejected reply from the Map-Server.
- o Reject: The reliable transport registration for the local EID database prefix was rejected by the Map-Server. From the ETR's perspective, no registration is due to the peer AND the peer is known to have rejected the registration.

The following events drive the state transitions:

- o DB creation: The local EID database entry for the EID prefix is created.

- o DB deletion: The local EID database entry for the EID prefix is deleted.
- o DB change: The mapping contents or authentication information for the local EID database entry changes.
- o Session up: The reliable transport session to the Map-Server is established.
- o Session down: The reliable transport session the Map-Server goes down.
- o Recv Refresh: A Registration refresh message is received from the Map-Server.
- o Recv ACK: A Registration Acknowledge message is received from the Map-Server.
- o Recv Rejected: A Registration Rejected message is received from the Map-Server.
- o Periodic timer: The timer that drives generation of periodic UDP Map-Register messages fires.

The state machine is:

Event	Prev State	
	No state	Periodic
DB creation [session down]	-> Periodic A1	N/A
DB creation [session up]	-> AckWait A2	N/A
DB deletion	N/A	-> No state A3
DB change	N/A	- A1
Session up	-	-> Stable A4
Session down	-	N/A
Recv Refresh	-	N/A
Recv Refresh [rejected]	-	N/A
Recv ACK	-	N/A
Recv Rejection	-	N/A
Timer	N/A	- A5

xTR per EID prefix per MS state machine

Event	Prev State		
	Stable	AckWait	Rejected
DB creation	N/A	N/A	N/A
DB deletion	-> No state A6	-> No state A6	-> No state
DB change	-> AckWait A2	- A2	-> AckWait A2
Session up	N/A	N/A	N/A
Session down	-> Periodic A7	-> Periodic A7	-> Periodic A7
Recv Refresh	-> AckWait A2	- A2	-> AckWait A2
Recv Refresh [rejected]	-	- A2	-> AckWait A2
Recv ACK	-	-> Stable	-> AckWait A2
Recv Rejection	-> Rejected	-> Rejected	-
Timer	N/A	N/A	N/A

xTR per EID prefix per MS state machine

Action descriptions:

- o A1: Start periodic registration timer with zero delay.
- o A2: Send Registration over reliable transport session.
- o A3: Send UDP registration with zero TTL.
- o A4: Stop periodic registration timer.

- o A7: Send UDP registration and start periodic registration timer with registration period.
- o A6: Send Registration with TTL zero over reliable transport session.
- o A7: Start periodic registration timer with registration period.

All timer start actions must be jittered.

When the reliable transport session is established the state machine moves into the Stable state without first registering the EID prefix over the reliable transport session. The subsequent refresh issued by the Map-Server will trigger the registration message to be sent. This model will allow future optimisations where the Map-Server may retain registration state from a previous instantiation of the reliable transport session with the ETR and only request the refresh of EID prefix state beyond some negotiated session progress marker.

Aa Map-Server authentication key change is treated as a DB change event and will result in triggering a new Registration message to be transmitted.

6.3. Map-Server Behavior

Received registrations create/update or delete mapping state.

A refresh for an unspecified prefix is sent when a session is first established to obtain the complete database contents from the ETR.

Refresh for rejected registrations sent (R bit set) when a new EID prefix is configured on the Map-Server.

Rejection sent to the ETR when an EID prefix that is registered is deconfigured.

Rejected Refresh (R bit set) sent when authentication for an EID prefix changes followed by a Rejection for existing registrations which fail authentication following change.

Mapping Notification message sent whenever the mapping for a registered or more specific prefix for which notifications are requested changes. ETR acknowledgement or rejection messaging for Mapping Notification is not required because the ETR decides how to process the message based on the registered mapping information. If the mapping information changes the resulting registration will trigger a new Mapping Notification message from the Map-Server.

7. Security Considerations

The LISP reliable transport session SHOULD be authenticated. On controlled RLOC networks that can guarantee that the source RLOC address of data packets cannot be spoofed, the authentication check can be a source address validation on the reliable transport packets. When the RLOC network does not provide such guarantees, reliable transport authentication SHOULD be used. Implementations SHOULD support the TCP Authentication Option (TCP-AO) [RFC5925] and SCTP Authenticated Chunks [RFC4895].

8. IANA Considerations

8.1. LISP Reliable Transport Message Types

Assignment of new LISP reliable transport message types is done according to the "IETF Review" model defined in [RFC5266].

The initial content of the registry should be as follows.

Type	Name	Reference
0-15	Reserved	This document
16	Error Notification	This document
17	Registration Message	This document
18	Registration Acknowledgement Message	This document
19	Registration Rejected Message	This document
20	Registration Refresh Message	This document
21	Mapping Notification Message	This document
22-30	Reserved for EID membership distribution	TBD
31-64999	Unassigned	
65000-65535	Reserved for Experimental Use	

8.2. Transport Protocol Port Numbers

TCP port 4342 already reserved for LISP CONS that is now obsolete. Repurpose for reliable transport over TCP. Reserve an SCTP port.

9. Acknowledgments

The authors would like to thank Noel Chiappa, Dino Farinacci, Jesper Skriver, Johnson Leong, Andre Pelletier and Les Ginsberg for their contributions to this document.

10. Normative References

- [I-D.ietf-lisp-lcaf]
Farinacci, D., Meyer, D., and J. Snijders, "LISP Canonical Address Format (LCAF)", draft-ietf-lisp-lcaf-10 (work in progress), June 2015.
- [I-D.kouvelas-lisp-reliable-transport]
Cassar, C., Kouvelas, I., and D. Lewis, "LISP Reliable Transport", draft-kouvelas-lisp-reliable-transport-02 (work in progress), March 2015.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC5266] Devarapalli, V. and P. Eronen, "Secure Connectivity and Mobility Using Mobile IPv4 and IKEv2 Mobility and Multihoming (MOBIKE)", BCP 136, RFC 5266, June 2008.
- [RFC6830] Farinacci, D., Fuller, V., Meyer, D., and D. Lewis, "The Locator/ID Separation Protocol (LISP)", RFC 6830, January 2013.
- [RFC6833] Fuller, V. and D. Farinacci, "Locator/ID Separation Protocol (LISP) Map-Server Interface", RFC 6833, January 2013.

Authors' Addresses

Chris Cassar
Cisco Systems
10 New Square Park
Bedfont Lakes, Feltham TW14 8HA
United Kingdom

Email: ccassar@cisco.com

Isidor Kouvelas
Arista Networks Inc.
5453 Great America Parkway
Santa Clara, CA 95054
USA

Email: kouvelas@cisco.com

Darrel Lewis
Cisco Systems
Tasman Drive
San Jose, CA 95134
USA

Email: darlewis@cisco.com

Jesus Arango
Cisco Systems
Tasman Drive
San Jose, CA 95134
USA

Email: jearango@cisco.com

Johnson Leong
Cisco Systems
Tasman Drive
San Jose, CA 95134
USA

Email: joleong@cisco.com