

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: August 25, 2022

D. Farinacci
lispers.net
C. Cantrell
Nexus
February 21, 2022

A Decent LISP Mapping System (LISP-Decent)
draft-farinacci-lisp-decent-09

Abstract

This draft describes how the LISP mapping system designed to be distributed for scale can also be decentralized for management and trust.

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 <https://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 August 25, 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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. Definition of Terms	3
3. Overview	4
4. Push-Based Mapping System	5
4.1. Components of a Pushed-Based LISP-Decent xTR	5
4.2. No LISP Protocol Changes	6
4.3. Configuration and Authentication	7
4.4. Core Seed-Group	7
5. Pull-Based Mapping System	9
5.1. Components of a Pulled-Based LISP-Decent xTR	9
5.2. Deployment Example	10
5.3. Management Considerations	11
6. Security Considerations	11
7. IANA Considerations	12
8. References	12
8.1. Normative References	12
8.2. Informative References	13
Appendix A. Acknowledgments	13
Appendix B. Document Change Log	14
B.1. Changes to draft-farinacci-lisp-decent-09	14
B.2. Changes to draft-farinacci-lisp-decent-08	14
B.3. Changes to draft-farinacci-lisp-decent-07	14
B.4. Changes to draft-farinacci-lisp-decent-06	14
B.5. Changes to draft-farinacci-lisp-decent-05	14
B.6. Changes to draft-farinacci-lisp-decent-04	14
B.7. Changes to draft-farinacci-lisp-decent-03	14
B.8. Changes to draft-farinacci-lisp-decent-02	15
B.9. Changes to draft-farinacci-lisp-decent-01	15
B.10. Changes to draft-farinacci-lisp-decent-00	15
Authors' Addresses	15

1. Introduction

The LISP architecture and protocols [RFC6830] introduces two new numbering spaces, Endpoint Identifiers (EIDs) and Routing Locators (RLOCs) which is intended to provide overlay network functionality. To map from EID to a set of RLOCs, a control-plane mapping system are used [RFC6836] [RFC8111]. These mapping systems are distributed in nature in their deployment for scalability but are centrally managed by a third-party entity, namely a Mapping System Provider (MSP). The entities that use the mapping system, such as data-plane xTRs, depend on and trust the MSP. They do not participate in the mapping system other than to register and retrieve information to/from the mapping system [RFC6833].

This document introduces a Decentralized Mapping System (DMS) so the xTRs can participate in the mapping system as well as use it. They can trust each other rather than rely on third-party infrastructure. The xTRs act as Map-Servers to maintain distributed state for scale and reducing attack surface.

2. Definition of Terms

Mapping System Provider (MSP): is an infrastructure service that deploys LISP Map-Resolvers and Map-Servers [RFC6833] and possibly ALT-nodes [RFC6836] or DDT-nodes [RFC8111]. The MSP can be managed by a separate organization other than the one that manages xTRs. This model provides a business separation between who manages and is responsible for the control-plane versus who manages the data-plane overlay service.

Decentralized Mapping System (DMS): is a mapping system entity that is not third-party to the xTR nodes that use it. The xTRs themselves are part of the mapping system. The state of the mapping system is fully distributed, decentralized, and the trust relies on the xTRs that use and participate in their own mapping system.

Pull-Based Mapping System: the mapping system is pull-based meaning that xTRs will lookup and register mappings by algorithmic transformation to locate which Map-Resolvers and Map-Servers are used. It is required that the lookup and registration uses a consistent algorithmic transformation function. Map-Registers are pushed to specific Map-Servers. Map-Requests are external lookups to Map-Resolvers on xTRs that do not participate in the mapping system and internal lookups when they do.

Modulus Value: this value is used in the Pull-Based Mapping System. It defines the number of map-server sets used for the mapping system. The modulus value is used to produce a Name Index used for a DNS lookup.

Name Index: this index value <index> is used in the Pull-Based Mapping System. For a mapping system that is configured with a map-server set of DNS names in the form of <name>.domain.com, the name index is prepended to <name> to form the lookup name <index>.<name>.domain.com. If the Modulus Value is 8, then the name indexes are 0 through 7.

Hash Mask: The Hash Mask is used in the Pull-Based Mapping System. It is a mask value with 1 bits left justified. The mask is used to select what high-order bits of an EID-prefix is used in the hash function.

Push-Based Mapping System: the mapping system is push-based meaning that xTRs will push registrations via IP multicast to a group of Map-Servers and do local lookups acting as their own Map-Resolvers.

Replication List Entry (RLE): is an RLOC-record format that contains a list of RLOCs that an ITR replicates multicast packets on a multicast overlay. The RLE format is specified in [RFC8060]. RLEs are used with the Pushed-Based mapping system.

Group Address EID: is an EID-record format that contains IPv4 (0.0.0.0/0, G) or IPv6 (0::/0, G) state. This state is encoded as a Multicast Info Type LCAF specified in [RFC8060]. Members of a seed-group send Map-Registers for (0.0.0.0/0, G) or (0::/0, G) with an RLOC-record that RLE encodes its RLOC address. Details are specified in [RFC8378].

Seed-Group: is a set of Map-Servers joined to a multicast group for the Push-Based Mapping system or are mapped by DNS names in a Pull-Based Mapping System. A core seed-group is used to bootstrap a set of LISP-Decent xTRs so they can learn about each other and use each other's mapping system service. A seed-group can be pull-based to bootstrap a push-based mapping system. That is, a set of DNS mapped map-servers can be used to join the mapping system's IP multicast group.

3. Overview

The clients of the Decentralized Mapping System (DMS) are also the providers of mapping state. Clients are typically ETRs that Map-Register EID-to-RLOC mapping state to the mapping database system. ITRs are clients in that they send Map-Requests to the mapping database system to obtain EID-to-RLOC mappings that are cached for data-plane use. When xTRs participate in a DMS, they are also acting as Map-Resolvers and Map-Servers using the protocol machinery defined in LISP control-plane specifications [RFC6833], [I-D.ietf-lisp-sec], and [I-D.ietf-lisp-ecdsa-auth]. The xTRs are not required to run the database mapping transport system protocols specified in [RFC6836] or [RFC8111].

This document will describe two decentralized and distributed mapping system mechanisms. A Push-Based Mapping System uses IP multicast so xTRs can find each other by locally joining an IP multicast group. A Pull-Based Mapping System uses DNS with an algorithmic transformation function so xTRs can find each other.

4. Push-Based Mapping System

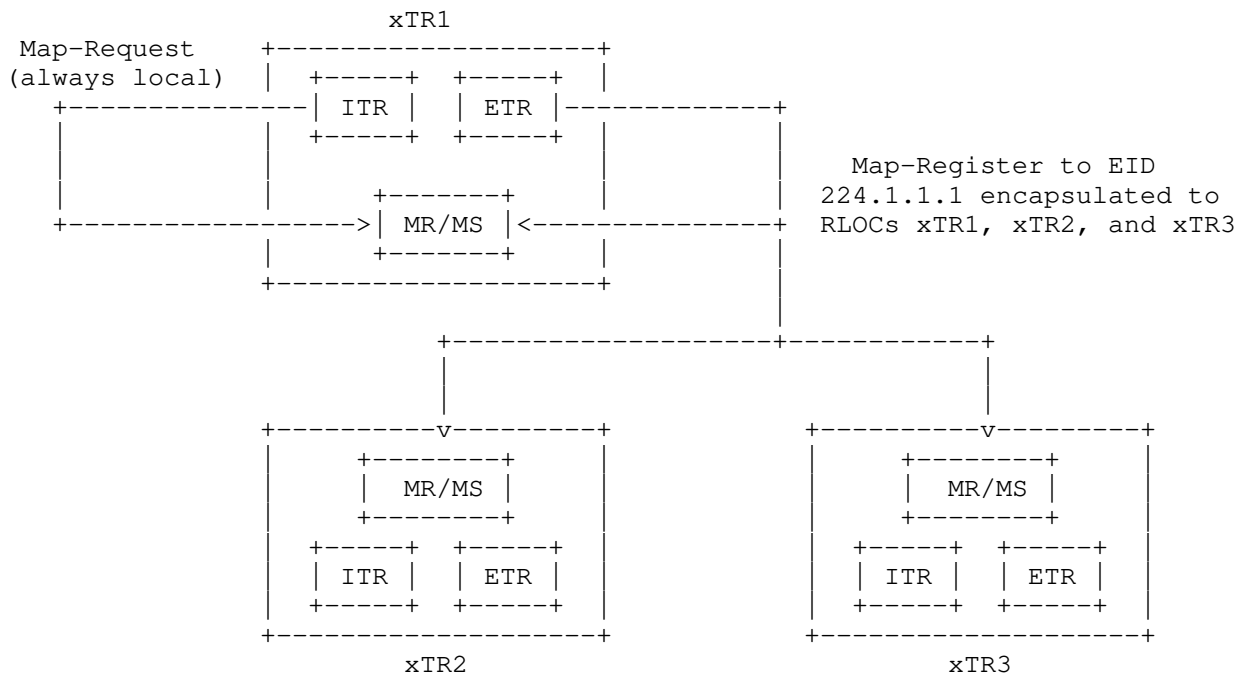
The xTRs are organized in a mapping-system group. The group is identified by an IPv4 or IPv6 multicast group address or using a pull-based approach in described in Section 5. When using multicast, the xTRs join the same multicast group and receive LISP control-plane messages addressed to the group. Messages sent to the multicast group are distributed when the underlay network supports IP multicast [RFC6831] or is achieved with the overlay multicast mechanism described in [RFC8378]. When overlay multicast is used and LISP Map-Register messages are sent to the group, they are LISP data encapsulated with a instance-ID set to 0xffffffff in the LISP header. The inner header of the encapsulated packet has the destination address set to the multicast group address and the outer header that is prepended has the destination address set to the RLOC of mapping system member. The members of the mapping system group are kept in the LISP data-plane map-cache so packets for the group can be replicated to each member RLOC.

All xTRs in a mapping system group will store the same registered mappings and maintain the state as Map-Servers normally do. The members are not only receivers of the multicast group but also send packets to the group.

4.1. Components of a Pushed-Based LISP-Decent xTR

When an xTR is configured to be a LISP-Decent xTR (or PxTR [RFC6832]), it runs the ITR, ETR, Map-Resolver, and Map-Server LISP network functions.

The following diagram shows 3 LISP-Decent xTRs joined to mapping system group 224.1.1.1. When the ETR function of xTR1 originates a Map-Register, it is sent to all xTRs (including itself) synchronizing all 3 Map-Servers in xTR1, xTR2, and xTR3. The ITR function can populate its map-cache by sending a Map-Request locally to its Map-Resolver so it can replicate packets to each RLOC for EID 224.1.1.1.



Note if any external xTR would like to use a Map-Resolver from the mapping system group, it only needs to have one of the LISP-Decent Map-Resolvers configured. By doing a looking to this Map-Resolver for EID 224.1.1.1, the external xTR could get the complete list of members for the mapping system group.

For future study, an external xTR could multicast the Map-Request to 224.1.1.1 and either one of the LISP-Decent Map-Resolvers would return a Map-Reply or the external xTR is prepared to receive multiple Map-Replies.

4.2. No LISP Protocol Changes

There are no LISP protocol changes required to support the push-based LISP-Decent set of procedures. However, an implementation that sends Map-Register messages to a multicast group versus a specific Map-Server unicast address must change to call the data-plane component so the ITR functionality in the node can encapsulate the Map-Register as a unicast packet to each member of the mapping system group.

An ITR SHOULD lookup its mapping system group address periodically to determine if the membership has changed. The ITR can also use the

pubsub capability documented in [I-D.ietf-lisp-pubsub] to be notified when a new member joins or leaves the multicast group.

4.3. Configuration and Authentication

When xTRs are joined to a multicast group, they must have their site registration configuration consistent. Any policy or authentication key material must be configured correctly and consistently among all members. When [I-D.ietf-lisp-ecdsa-auth] is used to sign Map-Register messages, public-keys can be registered to the mapping system group using the site authentication key mentioned above or using a different authentication key from the one used for registering EID records.

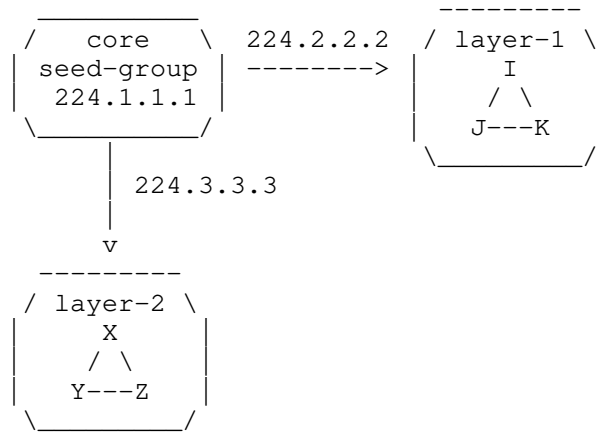
4.4. Core Seed-Group

A core seed-group can be discovered using a multicast group in a push-based system or a Map-Server set of DNS names in a pull-based system (see Section 5 for details).

When using multicast for the mapping system group, a core seed-group multicast group address can be preconfigured to bootstrap the decentralized mapping system. The group address (or DNS name that maps to a group address) can be explicitly configured in a few xTRs to start building up the registrations. Then as other xTRs come online, they can add themselves to the core seed-group by joining the seed-group multicast group.

Alternatively or additionally, new xTRs can join a new mapping system multicast group to form another layer of a decentralized mapping system. The group address and members of this new layer seed-group would be registered to the core seed-group address and stored in the core seed-group mapping system. Note each mapping system layer could have a specific function or a specific circle of trust.

This multi-layer mapping system can be illustrated:



Configured in xTRs A, B, and C (they make up the core seed-group):
 224.1.1.1 -> RLE: A, B, C

core seed-group DMS, mapping state in A, B, and C:

224.2.2.2 -> RLE: I, J, K

224.3.3.3 -> RLE: X, Y, Z

layer-1 seed-group DMS (inter-continental), mapping state in I, J, K:

EID1 -> RLOCs: i(1), j(2)

...

EIDn -> RLOCs: i(n), j(n)

layer-2 seed-group DMS (intra-continental), mapping sate in X, Y, Z::

EIDa -> RLOCs: x(1), y(2)

...

EIDz -> RLOCs: x(n), y(n)

The core seed-group multicast address 224.1.1.1 is configured in xTRs A, B and C so when each of them send Map-Register messages, they would all be able to maintain synchronized mapping state. Any EID can be registered to this DMS but in this example, seed-group multicast group EIDs are being registered only to find other mapping system groups.

For example, lets say that xTR I boots up and it wants to find its other peers in its mapping system group 224.2.2.2. Group address 224.2.2.2 is configured so xTR I knows what group to join for its mapping system group. But xTR I needs a mapping system to register to, so the core seed-group is used and available to receive Map-

Registers. The other xTRs J and K in the mapping system group do the same so when any of I, J or K needs to register EIDs, they can now send their Map-Register messages to group 224.2.2.2. Examples of EIDs being register are EID1 through EIDn shown above.

When Map-Registers are sent to group 224.2.2.2, they are encapsulated by the LISP data-plane by looking up EID 224.2.2.2 in the core seed-group mapping system. For the map-cache entry to be populated for 224.2.2.2, the data-plane must send a Map-Request so the RLOCs I, J, and K are cached for replication. To use the core seed-group mapping system, the data-plane must know of at least one of the RLOCs A, B, and/or C.

5. Pull-Based Mapping System

5.1. Components of a Pulled-Based LISP-Decent xTR

When an xTR is configured to be a LISP-Decent xTR (or PxTR [RFC6832]), it runs the ITR, ETR, Map-Resolver, and Map-Server LISP network functions.

Unlike the Push-Based Mapping System, the xTRs do not need to be organized by joining a multicast group. In a Pull-Based Mapping System, a hash function over an EID is used to identify which xTR is used as the Map-Resolver and Map-Server. The Domain Name System (DNS) [RFC1034] [RFC1035] is used as a resource discovery mechanism.

The RLOC addresses of the xTRs will be A and AAAA records for DNS names that map algorithmically from the hash of the EID. A SHA-256 hash function [RFC6234] over the following ASCII formatted EID string is used:

```
[<iid>]<eid>/<ml>
[<iid>]<group>/<gml>-<source>/<sml>
```

Where <iid> is the instance-ID and <eid> is the EID of any EID-type defined in [RFC8060]. And then the Modulus Value <mv> is used to produce the Name Index <index> used to build the DNS lookup name:

```
eid = "[<iid>]<eid>/<ml>"
index = hash.sha_256(eid) MOD mv
```

The Hash Mask is used to select what bits are used in the SHA-256 hash function. This is required to support longest match lookups in the mapping system. The same map-server set needs to be selected when looking up a more-specific EID found in the Map-Request message with one that could match a less-specific EID-prefix registered and found in the Map-Register message. For example, if an EID-prefix

[0]240.0.1.0/24 is registered to the mapping system and EID
[0]240.0.1.1/32 is looked up to match the registered prefix, a Hash
Mask of 8 bytes can be used to AND both the /32 or /24 entries to
produce the same hash string bits of "[0]240.0".

For (*,G) and (S,G) multicast entries in the mapping system, the hash
strings are:

```
sg-eid = "<[iid]><group>/<gml>-<source>/<sml>"
index = hash.sha_256(sg-eid) MOD mv

starg-eid = "<[iid]><group>/<gml>-0.0.0.0/0"
index = hash.sha_256(starg-eid) MOD mv
```

The Hash Mask MUST include the string "<[iid]><group>" and not string
<source>. So when looking up [0](2.2.2.2, 224.1.1.1) that will match
a (*, 224.1.1.1/32), the hash string produced with a Hash Mask of 12
bytes is "[0]224.1.1.1".

When the <index> is computed from a unicast or multicast EID, the DNS
lookup name becomes:

```
<index>.map-server.domain.com
```

When an xTR does a DNS lookup on the lookup name, it will send Map-
Register messages to all A and AAAA records for EID registrations.
For Map-Request messages, xTRs MAY round robin EID lookup requests
among the A and AAAA records.

5.2. Deployment Example

Here is an example deployment of a pull-based model. Let's say 4
map-server sets are provisioned for the mapping system. Therefore 4
distinct DNS names are allocated and a Modulus Value 4 is used. Each
DNS name is allocated Name Index 0 through 3:

```
0.map-server.lispers.net
1.map-server.lispers.net
2.map-server.lispers.net
3.map-server.lispers.net
```

The A records for each name can be assigned as:

```
0.map-server.lispers.net:
  A <rloc1-att>
  A <rloc2-verizon>
1.map-server.lispers.net:
  A <rloc1-bt>
  A <rloc2-dt>
2.map-server.lispers.net:
  A <rloc1-cn>
  A <rloc2-kr>
3.map-server.lispers.net:
  A <rloc1-au>
  A <rloc2-nz>
```

When an xTR wants to register "[1000]fd::2222", it hashes the EID string to produce, for example, hash value 0x66. Using the modulus value 4 (0x67 & 0x3) produces index 0x3, so the DNS name 3.map-server.lispers.net is used and a Map-Register is sent to <rloc1-au> and <rloc2-nz>.

Note that the pull-based method can be used for a core seed-group for bootstrapping a push-based mapping system where multicast groups are registered.

5.3. Management Considerations

There are no LISP protocol changes required to support the pull-based LISP-Decent set of procedures. However, an implementation SHOULD do periodic DNS lookups to determine if A records have changed for a DNS entry.

When xTRs derive Map-Resolver and Map-Server names from the DNS, they need to use the same Modulus Value otherwise some xTRs will lookup EIDs to the wrong place they were registered.

The Modulus Value can be configured or pushed to the LISP-Decent xTRs. A future version of this document will describe a push mechanism so all xTRs use a consistent modulus value.

6. Security Considerations

Refer to the Security Considerations section of [I-D.ietf-lisp-rfc6833bis] for a complete list of security mechanisms as well as pointers to threat analysis drafts.

7. IANA Considerations

At this time there are no specific requests for IANA.

8. References

8.1. Normative References

- [RFC1034] Mockapetris, P., "Domain names - concepts and facilities", STD 13, RFC 1034, DOI 10.17487/RFC1034, November 1987, <<https://www.rfc-editor.org/info/rfc1034>>.
- [RFC1035] Mockapetris, P., "Domain names - implementation and specification", STD 13, RFC 1035, DOI 10.17487/RFC1035, November 1987, <<https://www.rfc-editor.org/info/rfc1035>>.
- [RFC6234] Eastlake 3rd, D. and T. Hansen, "US Secure Hash Algorithms (SHA and SHA-based HMAC and HKDF)", RFC 6234, DOI 10.17487/RFC6234, May 2011, <<https://www.rfc-editor.org/info/rfc6234>>.
- [RFC6830] Farinacci, D., Fuller, V., Meyer, D., and D. Lewis, "The Locator/ID Separation Protocol (LISP)", RFC 6830, DOI 10.17487/RFC6830, January 2013, <<https://www.rfc-editor.org/info/rfc6830>>.
- [RFC6831] Farinacci, D., Meyer, D., Zwiebel, J., and S. Venaas, "The Locator/ID Separation Protocol (LISP) for Multicast Environments", RFC 6831, DOI 10.17487/RFC6831, January 2013, <<https://www.rfc-editor.org/info/rfc6831>>.
- [RFC6832] Lewis, D., Meyer, D., Farinacci, D., and V. Fuller, "Interworking between Locator/ID Separation Protocol (LISP) and Non-LISP Sites", RFC 6832, DOI 10.17487/RFC6832, January 2013, <<https://www.rfc-editor.org/info/rfc6832>>.
- [RFC6833] Fuller, V. and D. Farinacci, "Locator/ID Separation Protocol (LISP) Map-Server Interface", RFC 6833, DOI 10.17487/RFC6833, January 2013, <<https://www.rfc-editor.org/info/rfc6833>>.
- [RFC6836] Fuller, V., Farinacci, D., Meyer, D., and D. Lewis, "Locator/ID Separation Protocol Alternative Logical Topology (LISP+ALT)", RFC 6836, DOI 10.17487/RFC6836, January 2013, <<https://www.rfc-editor.org/info/rfc6836>>.

- [RFC8060] Farinacci, D., Meyer, D., and J. Snijders, "LISP Canonical Address Format (LCAF)", RFC 8060, DOI 10.17487/RFC8060, February 2017, <<https://www.rfc-editor.org/info/rfc8060>>.
- [RFC8111] Fuller, V., Lewis, D., Ermagan, V., Jain, A., and A. Smirnov, "Locator/ID Separation Protocol Delegated Database Tree (LISP-DDT)", RFC 8111, DOI 10.17487/RFC8111, May 2017, <<https://www.rfc-editor.org/info/rfc8111>>.
- [RFC8378] Moreno, V. and D. Farinacci, "Signal-Free Locator/ID Separation Protocol (LISP) Multicast", RFC 8378, DOI 10.17487/RFC8378, May 2018, <<https://www.rfc-editor.org/info/rfc8378>>.

8.2. Informative References

- [I-D.ietf-lisp-ecdsa-auth] Farinacci, D. and E. Nordmark, "LISP Control-Plane ECDSA Authentication and Authorization", draft-ietf-lisp-ecdsa-auth-06 (work in progress), August 2021.
- [I-D.ietf-lisp-pubsub] Rodriguez-Natal, A., Ermagan, V., Cabellos, A., Barkai, S., and M. Boucadair, "Publish/Subscribe Functionality for LISP", draft-ietf-lisp-pubsub-09 (work in progress), June 2021.
- [I-D.ietf-lisp-rfc6833bis] Farinacci, D., Maino, F., Fuller, V., and A. Cabellos, "Locator/ID Separation Protocol (LISP) Control-Plane", draft-ietf-lisp-rfc6833bis-30 (work in progress), November 2020.
- [I-D.ietf-lisp-sec] Maino, F., Ermagan, V., Cabellos, A., and D. Saucez, "LISP-Security (LISP-SEC)", draft-ietf-lisp-sec-25 (work in progress), December 2021.

Appendix A. Acknowledgments

The authors would like to thank the LISP WG for their review and acceptance of this draft.

The authors would also like to give a special thanks to Roman Shaposhnik for several discussions that occurred before the first draft was published.

Appendix B. Document Change Log

[RFC Editor: Please delete this section on publication as RFC.]

B.1. Changes to draft-farinacci-lisp-decent-09

- o Posted February 2022.
- o Update references and document expiry timer.

B.2. Changes to draft-farinacci-lisp-decent-08

- o Posted August 2021.
- o Update references and document expiry timer.

B.3. Changes to draft-farinacci-lisp-decent-07

- o Posted March 2021.
- o Update references and document expiry timer.

B.4. Changes to draft-farinacci-lisp-decent-06

- o Posted September 2020.
- o Update references and document expiry timer.

B.5. Changes to draft-farinacci-lisp-decent-05

- o Posted March 2020.
- o Update references and document expiry timer.

B.6. Changes to draft-farinacci-lisp-decent-04

- o Posted September 2019.
- o Update references and document expiry timer.

B.7. Changes to draft-farinacci-lisp-decent-03

- o Posted March 2019.
- o Introduce the Hash Mask which is used to grab common bits from a registered prefix and a lookup prefix.

- o Spec how multicast lookups are done in the pull-based mapping system.
- o Indicate the hash string includes the unicast EID mask-length and multicast group and source mask-lengths.

B.8. Changes to draft-farinacci-lisp-decent-02

- o Posted November 2018.
- o Changed references from peer-group to seed-group to make the algorithms in this document more like how blockchain networks initialize the peer-to-peer network.
- o Added pull mechanism to compliment the push mechanism. The pull mechanism could be used as a seed-group to bootstrap the push mechanism.

B.9. Changes to draft-farinacci-lisp-decent-01

- o Posted July 2018.
- o Document timer and reference update.

B.10. Changes to draft-farinacci-lisp-decent-00

- o Initial draft posted January 2018.

Authors' Addresses

Dino Farinacci
lispers.net
San Jose, CA
USA

Email: farinacci@gmail.com

Colin Cantrell
Nexus
Tempe, AZ
USA

Email: colin@nexus.io

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: May 8, 2019

N. Barry
Stellar Development Foundation
G. Losa
UCLA
D. Mazieres
Stanford University
J. McCaleb
Stellar Development Foundation
S. Polu
Stripe Inc.
November 4, 2018

The Stellar Consensus Protocol (SCP)
draft-mazieres-dinrg-scp-05

Abstract

SCP is an open Byzantine agreement protocol resistant to Sybil attacks. It allows Internet infrastructure stakeholders to reach agreement on a series of values without unanimous agreement on what constitutes the set of important stakeholders. A big differentiator from other Byzantine agreement protocols is that, in SCP, nodes determine the composition of quorums in a decentralized way: each node selects sets of nodes it considers large or important enough to speak for the whole network, and a quorum must contain such a set for each of its members.

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 <https://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 May 8, 2019.

Copyright Notice

Copyright (c) 2018 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 (<https://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. The Model	3
2.1. Slice infrastructures	3
2.2. Input and output	4
3. Protocol	5
3.1. Federated voting	5
3.2. Basic types	8
3.3. Quorum slices	9
3.4. Nominate message	10
3.5. Ballots	12
3.6. Prepare message	13
3.7. Commit message	17
3.8. Externalize message	18
3.9. Summary of phases	19
3.10. Message envelopes	20
4. Security considerations	21
5. Acknowledgments	21
6. References	22
6.1. Normative References	22
6.2. Informative References	22
Authors' Addresses	23

1. Introduction

Various aspects of Internet infrastructure depend on irreversible and transparent updates to data sets such as authenticated mappings [I-D.watson-dinrg-delmap]. Examples include public key certificates and revocations, transparency logs [RFC6962], preload lists for HSTS [RFC6797] and HPKP [RFC7469], and IP address delegation [I-D.paillisse-sidrops-blockchain].

The Stellar Consensus Protocol (SCP) specified in this draft allows Internet infrastructure stakeholders to collaborate in applying irreversible transactions to public state. SCP is an open Byzantine agreement protocol that resists Sybil attacks by allowing individual parties to specify minimum quorum memberships in terms of specific trusted peers. Each participant chooses combinations of peers on which to depend such that these combinations can be trusted in aggregate. The protocol guarantees safety so long as these dependency sets transitively overlap and contain sufficiently many honest nodes correctly obeying the protocol.

Though bad configurations are theoretically possible, several analogies provide an intuition for why transitive dependencies overlap in practice. For example, given multiple entirely disjoint Internet-protocol networks, people would have no trouble agreeing on the fact that the network containing the world's top web sites is the Internet. Such a consensus can hold even without unanimous agreement on what constitute the world's top web sites. Similarly, if network operators listed all the ASes from whom they would consider peering or transit worthwhile, the transitive closures of these sets would contain significant overlap, even without unanimous agreement on the "tier-1 ISP" designation. Finally, while different browsers and operating systems have slightly different lists of valid certificate authorities, there is significant overlap in the sets, so that a hypothetical system requiring validation from "all CAs" would be unlikely to diverge.

A more detailed abstract description of SCP and its rationale, including an English-language proof of safety, is available in [SCP]. In particular, that reference shows that a necessary property for safety, termed quorum intersection despite ill-behaved nodes, is sufficient to guarantee safety under SCP, making SCP optimally safe against Byzantine node failure for any given configuration.

This document specifies the end-system logic and wire format of the messages in SCP.

2. The Model

This section describes the configuration and input/output values of the consensus protocol.

2.1. Slice infrastructures

The SCP protocol achieves consensus on what we call a slice infrastructure, defined by a set of nodes and and, for each node, a set of quorum slices that determine quorum membership in a

decentralized way. Each `_node_` in has a digital signature key and is named by the corresponding public key, which we term a "NodeID".

Each node choses one or more quorum slices, which are sets of nodes that all include the node itself. A quorum slice represents a large or important enough set of peers that the node selecting the quorum slice believes the slice collectively speaks for the whole network.

A `_quorum_` is a non-empty set of nodes containing at least one quorum slice of each of its members. For instance, suppose "v1" has the single quorum slice "{v1, v2, v3}", while each of "v2", "v3", and "v4" has the single quorum slice "{v2, v3, v4}". In this case, "{v2, v3, v4}" is a quorum because it contains a slice for each member. On the other hand "{v1, v2, v3}" is not a quorum, because it does not contain a quorum slice for "v2" or "v3". The smallest quorum including "v1" in this example is the set of all nodes "{v1, v2, v3, v4}".

Unlike traditional Byzantine agreement protocols, nodes in SCP only care about quorums to which they belong themselves (and hence that contain at least one of their quorum slices). Intuitively, this is what protects nodes from Sybil attacks. In the example above, if "v3" deviates from the protocol, maliciously inventing 96 Sybils "v5, v6, ..., v100", the honest nodes' quorums will all still include one another, ensuring that "v1", "v2", and "v4" continue to agree on output values.

Every message in the SCP protocol specifies the sender's quorum slices. Hence, by collecting messages, a node dynamically learns what constitutes a quorum and can decide when a particular message has been sent by a quorum to which it belongs. (Again, nodes do not care about quorums to which they do not belong themselves.)

2.2. Input and output

SCP produces a series of output `_values_` for consecutively numbered `_slots_`. At the start of a slot, higher-layer software on each node supplies a candidate input value. Nodes then exchange protocol messages to agree on one or a combination of nodes' input values as the slot's output value. After a pause to assemble new input values, the process repeats for the next slot, with a 5-second interval between slots.

A value typically encodes a set of actions to apply to a replicated state machine. During the pause between slots, nodes accumulate the next set of actions, amortizing the cost of consensus on one slot over arbitrarily many individual state machine operations.

In practice, only one or a small number of nodes' input values actually affect the output value for any given slot. As discussed in Section 3.4, which nodes' input values to use depends on a cryptographic hash of the slot number and node public keys. A node's chances of affecting the output value depend on how often it appears in other nodes' quorum slices.

From SCP's perspective, values are just opaque byte arrays whose interpretation is left to higher-layer software. However, SCP requires a `_validity_` function (to check whether a value is valid) and a `_combining_` function that reduces multiple candidate values into a single `_composite_` value. When nodes nominate multiple values for a slot, SCP nodes invoke this function to converge on a single composite value. By way of example, in an application where values consist of sets of transactions, the combining function could take the union of transaction sets. Alternatively, if values represent a timestamp and a set of transactions, the combining function might pair the highest nominated timestamp with the transaction set that has the highest hash value.

3. Protocol

The protocol consists of exchanging digitally-signed messages bound to nodes' quorum slices. The format of all messages is specified using XDR [RFC4506]. In addition to quorum slices, messages compactly convey votes on sets of conceptual statements. The core technique of voting with quorum slices is termed `_federated voting_`. We describe federated voting next, then detail protocol messages in the subsections that follow.

The protocol goes through four phases: NOMINATE, PREPARE, COMMIT, and EXTERNALIZE. The NOMINATE and PREPARE phases run concurrently (though NOMINATE's messages are sent earlier and it ends before PREPARE ends). The COMMIT and EXTERNALIZE phases are exclusive, with COMMIT occurring immediately after PREPARE and EXTERNALIZE immediately after COMMIT.

3.1. Federated voting

Federated voting is a process through which nodes `_confirm_` statements. Not every attempt at federated voting may succeed--an attempt to vote on some statement "a" may get stuck, with the result that nodes can confirm neither "a" nor its negation "!a". However, when a node succeeds in confirming a statement "a", federated voting guarantees two things:

1. No two well-behaved nodes will confirm contradictory statements in any configuration and failure scenario in which any protocol

can guarantee safety for the two nodes (i.e., quorum intersection for the two nodes holds despite ill-behaved nodes).

2. If a quorum "I" is guaranteed safety by #1 even when all nodes in "!I" are malicious, and one node in "I" confirms a statement "a", then eventually every member of "I" will also confirm "a".

Intuitively, these conditions are key to ensuring agreement among nodes as well as a weak form of liveness (the non-blocking property [building-blocks]) that is compatible with the FLP impossibility result [FLP].

As a node "v" collects signed copies of a federated voting message "m" from peers, two thresholds trigger state transitions in "v" depending on the message. We define these thresholds as follows:

- o _quorum threshold_: When every member of a quorum to which "v" belongs (including "v" itself) has issued message "m"
- o _blocking threshold_: When at least one member of each of "v"'s quorum slices (a set that does not necessarily include "v" itself) has issued message "m"

Each node "v" can send several types of message with respect to a statement "a" during federated voting:

- o _vote_ "a" states that "a" is a valid statement and constitutes a promise by "v" not to vote for any contradictory statement, such as "!a".
- o _accept_ "a" says that nodes may or may not come to agree on "a", but if they don't, then the system has experienced a catastrophic set of Byzantine failures to the point that no quorum containing "v" consists entirely of correct nodes. (Nonetheless, accepting "a" is not sufficient to act on it, as doing so could violate agreement, which is worse than merely getting stuck from lack of a correct quorum.)
- o _vote-or-accept_ "a" is the disjunction of the above two messages. A node implicitly sends such a message if it sends either _vote_ "a" or _accept_ "a". Where it is inconvenient and unnecessary to differentiate between _vote_ and _accept_, a node can explicitly send a _vote-or-accept_ message.
- o _confirm_ "a" indicates that _accept_ "a" has reached quorum threshold at the sender. This message is interpreted the same as _accept_ "a", but allows recipients to optimize their quorum

checks by ignoring the sender's quorum slices, as the sender asserts it has already checked them.

Figure 1 illustrates the federated voting process. A node "v" votes for a valid statement "a" that doesn't contradict statements in past `_vote_` or `_accept_` messages sent by "v". When the `_vote_` message reaches quorum threshold, the node accepts "a". In fact, "v" accepts "a" if the `_vote-or-accept_` message reaches quorum threshold, as some nodes may accept "a" without first voting for it. Specifically, a node that cannot vote for "a" because it has voted for "a"'s negation "!a" still accepts "a" when the message `_accept_` "a" reaches blocking threshold (meaning assertions about "!a" have no hope of reaching quorum threshold barring catastrophic Byzantine failure).

If and when the message `_accept_` "a" reaches quorum threshold, then "v" has confirmed "a" and the federated vote has succeeded. In effect, the `_accept_` messages constitute a second vote on the fact that the initial vote messages succeeded. Once "v" enters the confirmed state, it may issue a `_confirm_` "a" message to help other nodes confirm "a" more efficiently by pruning their quorum search at "v".

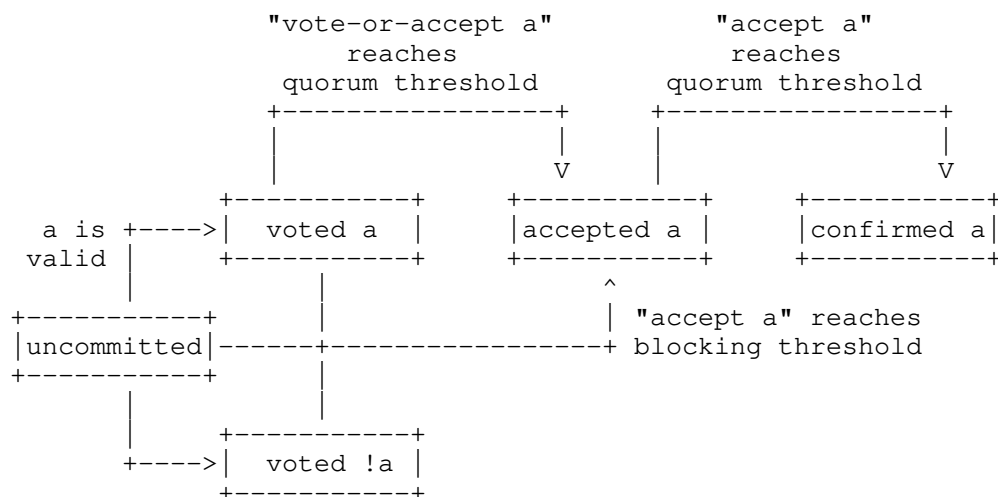


Figure 1: Federated voting process

Note several important invariants. A node may not vote for two contradictory statements or accept two contradictory statements. Moreover, a node may not vote for a statement that contradicts a message it has already accepted (which could lead to accepting a contradictory statement). However, a node is allowed to vote for one

statement and then accept a contradictory statement when a blocking threshold of accept messages contradicts the vote.

3.2. Basic types

SCP employs 32- and 64-bit integers, as defined below.

```
typedef unsigned int uint32;  
typedef int int32;  
typedef unsigned hyper uint64;  
typedef hyper int64;
```

SCP uses the SHA-256 cryptographic hash function [RFC6234], and represents hash values as a simple array of 32 bytes.

```
typedef opaque Hash[32];
```

SCP employs the Ed25519 digital signature algorithm [RFC8032]. For cryptographic agility, however, public keys are represented as a union type that can later be compatibly extended with other key types.

```
typedef opaque uint256[32];  
  
enum PublicKeyType  
{  
    PUBLIC_KEY_TYPE_ED25519 = 0  
};  
  
union PublicKey switch (PublicKeyType type)  
{  
    case PUBLIC_KEY_TYPE_ED25519:  
        uint256 ed25519;  
};  
  
// variable size as the size depends on the signature scheme used  
typedef opaque Signature<64>;
```

Nodes are public keys, while values are simply opaque arrays of bytes.

```
typedef PublicKey NodeID;  
  
typedef opaque Value<>;
```

3.3. Quorum slices

Theoretically a quorum slice can be an arbitrary set of nodes. However, arbitrary predicates on sets cannot be encoded concisely. Instead we specify quorum slices as any set of k-of-n members, where each of the n members can either be an individual node ID, or, recursively, another k-of-n set.

```
// supports things like: A,B,C, (D,E,F), (G,H, (I,J,K,L))
// only allows 2 levels of nesting
struct SCPSlices
{
    uint32 threshold;           // the k in k-of-n
    PublicKey validators<>;
    SCPSlices1 innerSets<>;
};
struct SCPSlices1
{
    uint32 threshold;           // the k in k-of-n
    PublicKey validators<>;
    SCPSlices2 innerSets<>;
};
struct SCPSlices2
{
    uint32 threshold;           // the k in k-of-n
    PublicKey validators<>;
};
```

Let "k" be the value of "threshold" and "n" the sum of the sizes of the "validators" and "innerSets" vectors in a message sent by some node "v". A message "m" sent by "v" reaches quorum threshold at "v" when three things hold:

1. "v" itself has issued (digitally signed) the message,
2. The number of nodes in "validators" who have signed "m" plus the number of "innerSets" that (recursively) meet this condition is at least "k", and
3. These three conditions apply (recursively) at some combination of nodes sufficient for condition #2.

A message reaches blocking threshold at "v" when the number of "validators" making the statement plus (recursively) the number "innerSets" reaching blocking threshold exceeds "n-k". (Blocking threshold depends only on the local node's quorum slices and hence does not require a recursive check on other nodes like step #3 above.)

As described in Section 3.10, every protocol message is paired with a cryptographic hash of the sender's "SCPSlices" and digitally signed. Inner protocol messages described in the next few sections should be understood to be received alongside such a quorum slice specification and digital signature.

3.4. Nominate message

For each slot, the SCP protocol begins in a NOMINATE phase, whose goal is to devise one or more candidate output values for the consensus protocol. In this phase, nodes send nomination messages comprising a monotonically growing set of values:

```
struct SCPNominate
{
    Value voted<>;      // X
    Value accepted<>;   // Y
};
```

The "voted" and "accepted" sets are disjoint; any value that is eligible for both sets is placed only in the "accepted" set.

"voted" consists of candidate values that the sender has voted to nominate. Each node progresses through a series of nomination rounds in which it may increase the set of values in its own "voted" field by adding the contents of the "voted" and "accepted" fields of "SCPNominate" messages received from a growing set of peers. In round "n" of slot "i", each node determines an additional peer whose nominated values it should incorporate in its own "SCPNominate" message as follows:

- o Let $G_i(m) = \text{SHA-256}(i \parallel m)$, where \parallel denotes the concatenation of serialized XDR values. Treat the output of G_i as a 256-bit binary number in big-endian format.
- o For each peer "v", define "weight(v)" as the fraction of quorum slices containing "v".
- o Define the set of nodes "neighbors(n)" as the set of nodes v for which $G_i(1 \parallel n \parallel v) < 2^{256} * \text{weight}(v)$, where "1" and "n" are both 32-bit XDR "int" values. Note that a node is always its own neighbor because conceptually a node belongs to all of its own quorum slices.
- o Define "priority(n, v)" as $G_i(2 \parallel n \parallel v)$, where "2" and "n" are both 32-bit XDR "int" values.

For each round "n" until nomination has finished (see below), a node starts echoing the available peer "v" with the highest value of "priority(n, v)" from among the nodes in "neighbors(n)". To echo "v", the node merges any valid values from "v"'s "voted" and "accepted" sets into its own "voted" set.

XXX - expand "voted" with only the 10 values with lowest Gi hash in any given round to avoid blowing out the message size?

Note that when echoing nominations, nodes must exclude and neither vote for nor accept values rejected by the higher-layer application's validity function. This validity function must not depend on state that can permanently differ across nodes. By way of example, it is okay to reject values that are syntactically ill-formed, that are semantically incompatible with the previous slot's value, that contain invalid digital signatures, that contain timestamps in the future, or that specify upgrades to unknown versions of the protocol. By contrast, the application cannot reject values that are incompatible with the results of a DNS query or some dynamically retrieved TLS certificate, as different nodes could see different results when doing such queries.

Nodes must not send an "SCPNominate" message until at least one of the "voted" or "accepted" fields is non-empty. When these fields are both empty, a node that has the highest priority among its neighbors in the current round (and hence should be echoing its own votes) adds the higher-layer software's input value to its "voted" field. Nodes that do not have the highest priority wait to hear "SCPNominate" messages from the nodes whose nominations they are echoing.

If a particular valid value "x" reaches quorum threshold in the messages sent by peers (meaning that every node in a quorum contains "x" either in the "voted" or the "accepted" field), then the node at which this happens moves "x" from its "voted" field to its "accepted" field and broadcasts a new "SCPNominate" message. Similarly, if "x" reaches blocking threshold in a node's peers' "accepted" field (meaning every one of a node's quorum slices contains at least one node with "x" in its "accepted" field), then the node adds "x" to its own "accepted" field (removing it from "voted" if applicable). These two cases correspond to the two conditions for entering the "accepted" state in Figure 1.

A node stops adding any new values to its "voted" set as soon as any value "x" reaches quorum threshold in the "accepted" fields of received "SCPNominate" messages. Following the terminology of Section 3.1, this condition corresponds to when the node confirms "x" as nominated. Note, however, that the node continues adding new values to "accepted" as appropriate. Doing so may lead to more

values becoming confirmed nominated even after the "voted" set is closed to new values.

A node always begins nomination in round "1". Round "n" lasts for "1+n" seconds, after which, if no value has been confirmed nominated, the node proceeds to round "n+1". A node continues to echo votes from the highest priority neighbor in prior rounds as well as the current round. In particular, until any value is confirmed nominated, a node continues expanding its "voted" field with values nominated by highest priority neighbors from prior rounds even when the values appeared after the end of those prior rounds.

As defined in the next two sections, the NOMINATE phase ends when a node has confirmed "prepare(b)" for some any ballot "b", as this is the point at which the nomination outcome no longer influences the protocol. Until this point, a node must continue to transmit "SCPNominate" messages as well as to expand its "accepted" set (even if "voted" is closed because some value has been confirmed nominated).

3.5. Ballots

Once there is a candidate on which to try to reach consensus, a node moves through three phases of balloting: PREPARE, COMMIT, and EXTERNALIZE. Balloting employs federated voting to chose between `_commit_` and `_abort_` statements for ballots. A ballot is a pair consisting of a counter and candidate value:

```
// Structure representing ballot <n, x>
struct SCPBallot
{
    uint32 counter; // n
    Value value;    // x
};
```

We use the notation "<n, x>" to represent a ballot with "counter == n" and "value == x".

Ballots are totally ordered with "counter" more significant than "value". Hence, we write "`b1 < b2`" to mean that either "`(b1.counter < b2.counter)`" or "`(b1.counter == b2.counter && b1.value < b2.value)`". Values are compared lexicographically as a strings of unsigned octets.

The protocol moves through federated voting on successively higher ballots until nodes confirm "`commit(b)`" for some ballot "b", at which point consensus terminates and outputs "b.value" for the slot. To ensure that only one value can be chosen for a slot and that the

protocol cannot get stuck if individual ballots get stuck, there are two restrictions on voting:

1. A node cannot vote for both "commit(b)" and "abort(b)" on the same ballot (the two outcomes are contradictory), and
2. A node may not vote for or accept "commit(b)" for any ballot "b" unless it has confirmed "abort" for every lesser ballot with a different value or already accepted "commit(b')" for some "b' < b" with "b'.value == b.value".

The second condition requires voting to abort large numbers of ballots before voting to commit a ballot "b". We call this `_preparing_` ballot "b", and introduce the following notation for the associated set of abort statements.

- o "prepare(b)" encodes an "abort" statement for every ballot less than "b" containing a value other than "b.value", i.e.,
`"prepare(b) = { abort(b1) | b1 < b AND b1.value != b.value }"`.
- o "vote prepare(b)" stands for a set of `_vote_` messages for every "abort" statement in "prepare(b)".
- o Similarly, "accept prepare(b)", "vote-or-accept prepare(b)", and "confirm prepare(b)" encode sets of `_accept_`, `_vote-or-accept_`, and `_confirm_` messages for every "abort" statement in "prepare(b)".

Using this terminology, a node must confirm "prepare(b)" before issuing a `_vote_` or `_accept_` message for the statement "commit(b)".

3.6. Prepare message

The first phase of balloting is the PREPARE phase. During this phase, as soon as a node has a valid candidate value (see the rules for "ballot.value" below), it begins sending the following message:

```
struct SCPPrepare
{
    SCPBallot ballot;           // current & highest prepare vote
    SCPBallot *prepared;       // highest accepted prepared ballot
    uint32 aCounter;           // lowest non-aborted ballot counter or 0
    uint32 hCounter;           // h.counter or 0 if h == NULL
    uint32 cCounter;           // c.counter or 0 if !c || !hCounter
};
```

This message compactly conveys the following (conceptual) federated voting messages:

- o "vote-or-accept prepare(ballot) "
- o If "prepared != NULL": "accept prepare(prepared) "
- o If "aCounter != 0": "accept abort(b) " for every "b" with "b.counter < aCounter"
- o If "hCounter != 0": "confirm prepare(<hCounter, ballot.value>)"
- o If "cCounter != 0": "vote commit(<n, ballot.value>)" for every "cCounter <= n <= hCounter"

Note that to be valid, an "SCPPPrepare" message must satisfy the following conditions:

- o If "prepared != NULL", then "prepared <= ballot" and "aCounter <= prepared.counter",
- o If "prepared == NULL", then "aCounter == 0", and
- o "cCounter <= hCounter <= ballot.counter".

Based on the federated vote messages received, each node keeps track of what ballots have been accepted and confirmed prepared. It uses these ballots to set the following fields of its own "SCPPPrepare" messages as follows.

ballot

The current ballot that a node is attempting to prepare and commit. The rules for setting each field are detailed below. Note that the "value" is updated when and only when "counter" changes.

ballot.counter

The counter is set according to the following rules:

- * Upon entering the PREPARE phase, the "counter" field is initialized to 1.
- * When a node sees messages from a quorum to which it belongs such that each message's "ballot.counter" is greater than or equal to the local "ballot.counter", the node arms a timer to fire in a number of seconds equal to its "ballot.counter + 1" (so the timeout lengthens linearly as the counter increases). Note that for the purposes of determining whether a quorum has a particular "ballot.counter", a node considers "ballot" fields in "SCPPPrepare" and "SCPCommit" messages. It also considers

"SCPExternalize" messages to convey an implicit "ballot.counter" of "infinity".

- * If the timer fires, a node increments the ballot counter by 1.
- * If nodes forming a blocking threshold all have "ballot.counter" values greater than the local "ballot.counter", then the local node immediately cancels any pending timer, increases "ballot.counter" to the lowest value such that this is no longer the case, and if appropriate according to the rules above arms a new timer. Note that the blocking threshold may include ballots from "SCPCommit" messages as well as "SCPExternalize" messages, which implicitly have an infinite ballot counter.
- * **Exception**: To avoid exhausting "ballot.counter", its value must always be less than 1,000 plus the number of seconds a node has been running SCP on the current slot. Should any of the above rules require increasing the counter beyond this value, a node either increases "ballot.counter" to the maximum permissible value, or, if it is already at this maximum, waits up to one second before increasing the value.

ballot.value

Each time the ballot counter is changed, the value is also recomputed as follows:

- * If any ballot has been confirmed prepared, then "ballot.value" is taken to be "h.value" for the highest confirmed prepared ballot "h". (Note that once this is the case, the node can stop sending "SCPNominate" messages, as "h.value" supersedes any output of the nomination protocol.)
- * Otherwise (if no such "h" exists), if one or more values are confirmed nominated, then "ballot.value" is taken as the output of the deterministic combining function applied to all confirmed nominated values. Note that because the NOMINATE and PREPARE phases run concurrently, the set of confirmed nominated values may continue to grow during balloting, changing "ballot.value" even if no ballots are confirmed prepared.
- * Otherwise, if no ballot is confirmed prepared and no value is confirmed nominated, but the node has accepted a ballot prepared (because "prepare(b)" meets blocking threshold for some ballot "b"), then "ballot.value" is taken as the value of the highest such accepted prepared ballot.

- * Otherwise, if no value is confirmed nominated and no value is accepted prepared, then a node cannot yet send an "SCPPprepare" message and must continue sending only "SCPNominate" messages.

prepared

The highest accepted prepared ballot not exceeding the "ballot" field, or NULL if no ballot has been accepted prepared. Recall that ballots with equal counters are totally ordered by the value. Hence, if "ballot = <n, x>" and the highest prepared ballot is "<n, y>" where "x < y", then the "prepared" field in sent messages must be set to "<n-1, y>" instead of "<n, y>", as the latter would exceed "ballot". In the event that "n = 1", the prepared field may be set to "<0, y>", meaning 0 is a valid "prepared.counter" even though it is not a valid "ballot.counter". It is possible to confirm "prepare(<0, y>)", in which case the next "ballot.value" is set to "y". However, it is not possible to vote to commit a ballot with counter 0.

aCounter

The lowest counter such that all ballots with lower counters have been accepted aborted. This value is set whenever "prepared.value" changes, since the definition of prepare implies that all ballots below the lesser of two prepared ballots have been aborted. Specifically, if the value of "prepared" just changed from "oldPrepared" where "prepared.value != oldPrepared.value", then "aCounter" is set to "oldPrepared.counter" if "oldPrepared.value < prepared.value", and "oldPrepared.counter+1" otherwise.

hCounter

If "h" is the highest confirmed prepared ballot and "h.value == ballot.value", then this field is set to "h.counter". Otherwise, if no ballot is confirmed prepared or if "h.value != ballot.value", then this field is 0. Note that by the rules above, if "h" exists, then "ballot.value" will be set to "h.value" the next time "ballot" is updated.

cCounter

The value "cCounter" is maintained based on an internally-maintained _commit ballot_ "c", initially "NULL". "cCounter" is 0 while "c == NULL" or "hCounter == 0", and is "c.counter" otherwise. "c" is updated as follows:

- * If either "(prepared > c && prepared.value != c.value)" or "(aCounter > c.counter)", then reset "c = NULL".
- * If "c == NULL" and "hCounter == ballot.counter" (meaning "ballot" is confirmed prepared), then set "c" to "ballot".

Note these rules preserve the invariant that a node cannot vote for contradictory statements (namely committing and aborting the same ballot) by conservatively assuming a node may have voted to abort anything below "ballot". Hence, whenever "c" changes, it can either change to "NULL" or to "ballot", but is never set to anything below the current "ballot".

A node leaves the PREPARE phase and proceeds to the COMMIT phase when there is some ballot "b" for which the node confirms "prepare(b)" and accepts "commit(b)". (If nodes never changed quorum slice mid-protocol, it would suffice to accept "commit(b)". Also waiting to confirm "prepare(b)" makes it easier to recover from liveness failures by removing Byzantine faulty nodes from quorum slices.)

3.7. Commit message

In the COMMIT phase, a node has accepted "commit(b)" for some ballot "b", and must confirm that statement to act on the value in "b.counter". A node sends the following message in this phase:

```
struct SCPCommit
{
    SCPBallot ballot;           // b
    uint32 preparedCounter;    // prepared.counter
    uint32 hCounter;           // h.counter
    uint32 cCounter;           // c.counter
};
```

The message conveys the following federated vote messages, where "infinity" is 2^{32} (a value greater than any ballot counter representable in serialized form):

- o "accept commit(<n, ballot.value>)" for every "cCounter <= n <= hCounter"
- o "vote-or-accept prepare(<infinity, ballot.value>)"
- o "accept prepare(<preparedCounter, ballot.value>)"
- o "confirm prepare(<hCounter, ballot.value>)"
- o "vote commit(<n, ballot.value>)" for every "n >= cCounter"

A node computes the fields in the "SCPCommit" messages it sends as follows:

ballot

This field is maintained identically to how it is maintained in the PREPARE phase, though "ballot.value" can no longer change, only "ballot.counter". Note that the value "ballot.counter" does not figure in any of the federated voting messages. The purpose of continuing to update and send this field is to assist other nodes still in the PREPARE phase in synchronizing their counters.

preparedCounter

This field is the counter of the highest accepted prepared ballot--maintained identically to the "prepared" field in the PREPARE phase. Since the "value" field will always be the same as "ballot", only the counter is sent in the COMMIT phase.

cCounter

The counter of the lowest ballot "c" for which the node has accepted "commit(c)". (No value is included in messages since "c.value == ballot.value".)

hCounter

The counter of the highest ballot "h" for which the node has accepted "commit(h)". (No value is included in messages since "h.value == ballot.value".)

As soon as a node confirms "commit(b)" for any ballot "b", it moves to the EXTERNALIZE phase.

3.8. Externalize message

A node enters the EXTERNALIZE phase when it confirms "commit(b)" for any ballot "b". As soon as this happens, SCP outputs "b.value" as the value of the current slot. In order to help other nodes achieve consensus on the slot more quickly, a node reaching this phase also sends the following message:

```
struct SCPExternalize
{
    SCPBallot commit;           // c
    uint32 hCounter;           // h.counter
};
```

An "SCPExternalize" message conveys the following federated voting messages:

- o "accept commit(<n, commit.value>)" for every "n >= commit.counter"
- o "confirm commit(<n, commit.value>)" for every "commit.counter <= n <= hCounter"

- o "accept prepare(<infinity, commit.value>)"
- o "confirm prepare(<hCounter, commit.value>)"

The fields are set as follows:

commit

The lowest confirmed committed ballot.

hCounter

The counter of the highest confirmed committed ballot.

3.9. Summary of phases

Table 1 summarizes the phases of SCP for each slot. The NOMINATE and PREPARE phases begin concurrently. However, a node initially does not send "SCPPrepate" messages but only listens for ballot messages in case "accept prepare(b)" reaches blocking threshold for some ballot "b". The COMMIT and EXTERNALIZE phases then run in turn after PREPARE ends. A node may externalize (act upon) a value as soon as it enters the EXTERNALIZE phase.

The point of "SCPEexternalize" messages is to help stragglers catch up more quickly. As such, the EXTERNALIZE phase never ends. Rather, a node should archive an "SCPEexternalize" message for as long as it retains slot state.

Phase	Begin	End
NOMINATE	previous slot externalized and 5 seconds have elapsed since NOMINATE ended for that slot	some ballot is confirmed prepared
PREPARE	begin with NOMINATE, but send "SCPPrepate" only once some value confirmed nominated or accept "prepare(b)" for some ballot b	accept "commit(b)" for some ballot "b"
COMMIT	accept "commit(b)" for some ballot "b"	confirm "commit(b)" for some ballot "b"
EXTERNALIZE	confirm "commit(b)" for some ballot "b"	slot state garbage-collected

Table 1: Phases of SCP for a slot

3.10. Message envelopes

In order to provide full context for each signed message, all signed messages are part of an "SCPStatement" union type that includes the "slotIndex" naming the slot to which the message applies, as well as the "type" of the message. A signed message and its signature are packed together in an "SCPEnvelope" structure.

```
enum SCPStatementType
{
    SCP_ST_PREPARE = 0,
    SCP_ST_COMMIT = 1,
    SCP_ST_EXTERNALIZE = 2,
    SCP_ST_NOMINATE = 3
};

struct SCPStatement
{
    NodeID nodeID;          // v (node signing message)
    uint64 slotIndex;       // i
    Hash quorumSetHash;     // hash of serialized SCPSlices

    union switch (SCPStatementType type)
    {
        case SCP_ST_PREPARE:
            SCPPrepare prepare;
        case SCP_ST_COMMIT:
            SCPCommit commit;
        case SCP_ST_EXTERNALIZE:
            SCPExternalize externalize;
        case SCP_ST_NOMINATE:
            SCPNominate nominate;
    }
    pledges;
};

struct SCPEnvelope
{
    SCPStatement statement;
    Signature signature;
};
```

4. Security considerations

If nodes do not pick quorum slices well, the protocol will not be safe.

5. Acknowledgments

The Stellar development foundation supported development of the protocol and produced the first production deployment of SCP. The IRTF DIN group including Dirk Kutscher, Sydney Li, Colin Man, Piers Powlesland, Melinda Shore, and Jean-Luc Watson helped with the framing and motivation for this specification. The mobilecoin team contributed the "aCounter" optimization. We also thank Bob

Glickstein for finding bugs in drafts of this document and offering many useful suggestions.

6. References

6.1. Normative References

- [RFC4506] Eisler, M., Ed., "XDR: External Data Representation Standard", STD 67, RFC 4506, DOI 10.17487/RFC4506, May 2006, <<https://www.rfc-editor.org/info/rfc4506>>.
- [RFC6234] Eastlake 3rd, D. and T. Hansen, "US Secure Hash Algorithms (SHA and SHA-based HMAC and HKDF)", RFC 6234, DOI 10.17487/RFC6234, May 2011, <<https://www.rfc-editor.org/info/rfc6234>>.
- [RFC8032] Josefsson, S. and I. Liusvaara, "Edwards-Curve Digital Signature Algorithm (EdDSA)", RFC 8032, DOI 10.17487/RFC8032, January 2017, <<https://www.rfc-editor.org/info/rfc8032>>.

6.2. Informative References

- [building-blocks] Song, Y., van Renesse, R., Schneider, F., and D. Dolev, "The Building Blocks of Consensus", 9th International Conference on Distributed Computing and Networking pp. 54-72, 2008.
- [FLP] Fischer, M., Lynch, N., and M. Lynch, "Impossibility of Distributed Consensus with One Faulty Process", Journal of the ACM 32(2):374-382, 1985.
- [I-D.paillisse-sidrops-blockchain] Paillisse, J., Rodriguez-Natal, A., Ermagan, V., Maino, F., Vegoda, L., and A. Cabellos-Aparicio, "An analysis of the applicability of blockchain to secure IP addresses allocation, delegation and bindings.", draft-paillisse-sidrops-blockchain-02 (work in progress), June 2018.
- [I-D.watson-dinrg-delmap] Watson, J., Li, S., and C. Man, "Delegated Distributed Mappings", draft-watson-dinrg-delmap-01 (work in progress), October 2018.

- [RFC6797] Hodges, J., Jackson, C., and A. Barth, "HTTP Strict Transport Security (HSTS)", RFC 6797, DOI 10.17487/RFC6797, November 2012, <<https://www.rfc-editor.org/info/rfc6797>>.
- [RFC6962] Laurie, B., Langley, A., and E. Kasper, "Certificate Transparency", RFC 6962, DOI 10.17487/RFC6962, June 2013, <<https://www.rfc-editor.org/info/rfc6962>>.
- [RFC7469] Evans, C., Palmer, C., and R. Sleevi, "Public Key Pinning Extension for HTTP", RFC 7469, DOI 10.17487/RFC7469, April 2015, <<https://www.rfc-editor.org/info/rfc7469>>.
- [SCP] Mazieres, D., "The Stellar Consensus Protocol: A Federated Model for Internet-level Consensus", Stellar Development Foundation whitepaper , 2015, <<https://www.stellar.org/papers/stellar-consensus-protocol.pdf>>.

Authors' Addresses

Nicolas Barry
Stellar Development Foundation
170 Capp St., Suite A
San Francisco, CA 94110
US

Email: nicolas@stellar.org

Giuliano Losa
UCLA
3753 Keystone Avenue #10
Los Angeles, CA 90034
US

Email: giuliano@cs.ucla.edu

David Mazieres
Stanford University
353 Serra Mall, Room 290
Stanford, CA 94305
US

Email: dm@uun.org

Jed McCaleb
Stellar Development Foundation
170 Capp St., Suite A
San Francisco, CA 94110
US

Email: jed@stellar.org

Stanislas Polu
Stripe Inc.
185 Berry Street, Suite 550
San Francisco, CA 94107
US

Email: stan@stripe.com