

Network Working Group
Internet-Draft
Updates: 6126bis (if approved)
Intended status: Standards Track
Expires: January 3, 2019

A. Decimo
IRIF, University of Paris-Diderot
D. Schinazi
Apple Inc.
J. Chroboczek
IRIF, University of Paris-Diderot
July 2, 2018

Babel Routing Protocol over Datagram Transport Layer Security
draft-decimo-babel-dtls-01

Abstract

This documents describes how to use Datagram Transport Layer Security (DTLS) to secure the Babel Routing Protocol.

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 January 3, 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
1.1. Specification of Requirements	3
2. Operation of the Protocol	3
3. Handling protected and unprotected data	3
3.1. Cleartext and DTLS on the same port	3
3.2. Cleartext and DTLS on separate ports	4
4. Establishing and handling Babel over DTLS sessions	4
4.1. Session Initiation	4
4.2. Transmission	4
4.3. Reception	5
4.4. Neighbour flush	5
5. Interface MTU Issues	5
6. IANA Considerations	6
7. Security Considerations	6
8. References	6
8.1. Normative References	6
8.2. Informative References	7
Appendix A. Performance Considerations	7
Authors' Addresses	8

1. Introduction

Babel over DTLS is a security protocol for the Babel routing protocol [RFC6126bis], that uses Datagram Transport Layer Security (DTLS) [RFC6347]. This document describes how to protect Babel with Babel over DTLS.

The motivation for proposing Babel over DTLS is that DTLS provides a sub-layer of security that is well-defined, whose security has been shown, and that has multiple implementations. Babel over DTLS has the following properties, inherited from DTLS:

- o authentication of peers;
- o integrity of data;
- o confidentiality of data;
- o use of asymmetric keys.

The main change to the Babel protocol is that Babel over DTLS requires most packets to be sent over unicast.

A malicious entity in range of a non-secured deployment of Babel can learn properties of the network, but also reroute legitimate traffic by advertising routes with a low metric.

1.1. Specification of Requirements

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. Operation of the Protocol

At first sight, there are incompatibilities between Babel and DTLS. Babel is a pure peer-to-peer protocol, while DTLS is a two parties client-server protocol. A Babel implementation typically uses unicast and multicast, while DTLS can only protect unicast.

The problem of assigning a client or server role for a DTLS handshake to Babel nodes is solved by a simple arbitrary choice. The addresses of the two nodes are compared, and the node with the lowest address acts as the server.

Babel is sufficiently flexible to work almost without multicast. In Babel over DTLS, almost all packets are sent via unicast. Packets that would have been sent via multicast need to be duplicated and sent to all of the original recipients via unicast. The cost of duplication is balanced by the fact that on networks with low throughput, unicast is often far more effective than multicast. Only neighbour discovery packets are sent via multicast, as they do not represent a security threat.

3. Handling protected and unprotected data

A Babel node needs to receive unprotected data for bootstrapping reasons, as well as protected data. Protected and unprotected traffic needs to be differentiated.

[NOTE TO READER: AUTHORS ARE CONSIDERING THESE TWO OPTIONS BUT ONLY ONE WILL BE RETAINED IN THE FINAL DOCUMENT.]

3.1. Cleartext and DTLS on the same port

In this approach, Babel and Babel over DTLS traffic is received on the same port. The DTLS client port, the DTLS server port, and the Babel port (6696) are equal. When a packet is received, it is unconditionally treated as a DTLS packet and decrypted.

- o If the decryption is successful, the decrypted content is parsed as a Babel packet and the node acts on it.

- o Otherwise, the packet is parsed as a Babel packet and the node MUST silently ignore all TLVs except Hello and IHU.

Since the source port is fixed as 6696, a node that loses its DTLS state (e.g. if it reboots), will reuse the same source and destination ports for the new session. In order to avoid discarding these new packets, nodes receiving an unexpected DTLS ClientHello MUST proceed with a new handshake and MUST NOT destroy the existing session until the new session's handshake completes to avoid denial of service attacks (Section 4.2.8 of [RFC6347]).

3.2. Cleartext and DTLS on separate ports

In this approach, a different port (number TBD) is allocated by IANA for Babel over DTLS traffic. The Babel over DTLS server listens on this port and Babel over DTLS clients use an ephemeral source port to initiate outbound DTLS connections. Unprotected Babel messages are sent and received over the standard Babel port (6696). When parsing unprotected packets, all Babel TLVs except Hello and IHU MUST be silently ignored.

4. Establishing and handling Babel over DTLS sessions

4.1. Session Initiation

When a node A acquires a new neighbour B (e.g. when A first receives a Babel packet from B, see Section 3.4 of [RFC6126bis]),

- o if the IP address A uses to send and receive Babel packets is smaller than the source IP address of the received Babel packet from B, A initialises its DTLS state as a server for peer B;
- o otherwise, A initialises its DTLS state as a client for peer B, and initiates a DTLS handshake.

Once the handshake succeeds and a DTLS session is established, nodes send all unicast Babel messages over DTLS.

4.2. Transmission

Since DTLS cannot secure multicast, nodes SHOULD send all TLVs over unicast DTLS, if possible. All TLVs that are not Hello nor IHU MUST be sent over unicast DTLS. Hello and IHU TLVs MAY be sent either over unicast DTLS or unprotected multicast. Nodes MUST NOT send any unprotected packets over unicast.

4.3. Reception

Packets received over unicast DTLS are parsed the same way as any packets in the original specification of Babel. Nodes MUST parse unprotected packets received over multicast, however they MUST silently ignore any TLV that is not Hello or IHU. Unprotected packets received over unicast MUST be silently ignored.

4.4. Neighbour flush

When a neighbour entry is flushed from the neighbour table (Appendix A of [RFC6126bis]), its associated DTLS state SHOULD be discarded. The node MAY send a DTLS close_notify alert to the neighbour.

5. Interface MTU Issues

Compared to normal Babel, DTLS adds at least 13 octets of header, plus cipher and authentication overhead to every packet. This reduces the size of the Babel payload that can be carried.

As stated in Section 4 of [RFC6126bis], in order to minimise the number of packets being sent while avoiding lower-layer fragmentation, a Babel node SHOULD attempt to maximise the size of the packets it sends, up to the outgoing interface's MTU adjusted for lower-layer headers (28 octets for UDP over IPv4, 48 octets for UDP over IPv6). It MUST NOT send packets larger than the attached interface's MTU adjusted for lower-layer headers or 512 octets, whichever is larger, but not exceeding $2^{16} - 1$ adjusted for lower-layer headers. Every Babel speaker MUST be able to receive packets that are as large as any attached interface's MTU adjusted for lower-layer headers or 512 octets, whichever is larger. Babel packets MUST NOT be sent in IPv6 Jumbograms.

These requirements are retained by this specification, but are extended to take DTLS overhead into account as follows. The Babel node MUST ensure that the DTLS datagram size does not exceed the interface MTU, i.e., each DTLS record MUST fit within a single datagram, as required by [RFC6347]. A Babel node MUST consider the amount of record expansion expected by the DTLS processing when calculating the maximum size of Babel packet that fits within the interface MTU. The overhead can be computed as DTLS overhead of 13 octets + authentication overhead of the negotiated DTLS cipher suite + block padding (Section 4.1.1.1 of [RFC6347]).

6. IANA Considerations

If the final version of this specification uses the standard Babel port for unprotected packets and DTLS Section 3.1, no actions are required from IANA.

If the final version of this specification uses separate ports for unprotected packets and DTLS Section 3.2, IANA is requested to assign a UDP port with label "Babel_DTLS".

7. Security Considerations

The interaction between two Babel peers requires Datagram Transport Layer Security (DTLS) with a cipher suite offering confidentiality protection. The guidance given in [RFC7525] MUST be followed to avoid attacks on DTLS. The DTLS client SHOULD use the TLS Certificate Status Request extension (Section 8 of [RFC6066]).

A malicious client might attempt to perform a high number of DTLS handshakes with a server. As the clients are not uniquely identified by the protocol and can be obfuscated with IPv4 address sharing and with IPv6 temporary addresses, a server needs to mitigate the impact of such an attack. Such mitigation might involve rate limiting handshakes from a given subnet or more advanced DoS/DDoS avoidance techniques beyond the scope of this document.

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC6126bis] Chroboczek, J. and D. Schinazi, "The Babel Routing Protocol", Internet Draft draft-ietf-babel-rfc6126bis-05, October 2017.
- [RFC6347] Rescorla, E. and N. Modadugu, "Datagram Transport Layer Security Version 1.2", RFC 6347, DOI 10.17487/RFC6347, January 2012, <<https://www.rfc-editor.org/info/rfc6347>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

8.2. Informative References

- [RFC6066] Eastlake 3rd, D., "Transport Layer Security (TLS) Extensions: Extension Definitions", RFC 6066, DOI 10.17487/RFC6066, January 2011, <<https://www.rfc-editor.org/info/rfc6066>>.
- [RFC7250] Wouters, P., Ed., Tschofenig, H., Ed., Gilmore, J., Weiler, S., and T. Kivinen, "Using Raw Public Keys in Transport Layer Security (TLS) and Datagram Transport Layer Security (DTLS)", RFC 7250, DOI 10.17487/RFC7250, June 2014, <<https://www.rfc-editor.org/info/rfc7250>>.
- [RFC7525] Sheffer, Y., Holz, R., and P. Saint-Andre, "Recommendations for Secure Use of Transport Layer Security (TLS) and Datagram Transport Layer Security (DTLS)", BCP 195, RFC 7525, DOI 10.17487/RFC7525, May 2015, <<https://www.rfc-editor.org/info/rfc7525>>.
- [RFC7918] Langley, A., Modadugu, N., and B. Moeller, "Transport Layer Security (TLS) False Start", RFC 7918, DOI 10.17487/RFC7918, August 2016, <<https://www.rfc-editor.org/info/rfc7918>>.
- [RFC7924] Santesson, S. and H. Tschofenig, "Transport Layer Security (TLS) Cached Information Extension", RFC 7924, DOI 10.17487/RFC7924, July 2016, <<https://www.rfc-editor.org/info/rfc7924>>.
- [RFC8094] Reddy, T., Wing, D., and P. Patil, "DNS over Datagram Transport Layer Security (DTLS)", RFC 8094, DOI 10.17487/RFC8094, February 2017, <<https://www.rfc-editor.org/info/rfc8094>>.

Appendix A. Performance Considerations

To reduce the number of octets taken by the DTLS handshake, especially the size of the certificate in the ServerHello (which can be several kilobytes), Babel peers can use raw public keys [RFC7250] or the Cached Information Extension [RFC7924]. The Cached Information Extension avoids transmitting the server's certificate and certificate chain if the client has cached that information from a previous TLS handshake. TLS False Start [RFC7918] can reduce round trips by allowing the TLS second flight of messages (ChangeCipherSpec) to also contain the (encrypted) Babel packet.

These performance considerations were inspired from the ones for DNS over DTLS [RFC8094].

Authors' Addresses

Antonin Decimo
IRIF, University of Paris-Diderot
Paris
France

Email: antonin.decimo@gmail.com

David Schinazi
Apple Inc.
One Apple Park Way
Cupertino, California 95014
USA

Email: dschinazi@apple.com

Juliusz Chroboczek
IRIF, University of Paris-Diderot
Case 7014
75205 Paris Cedex 13
France

Email: jch@irif.fr

Network Working Group
Internet-Draft

Obsoletes: 7298 (if approved)
Updates: 6126bis (if approved)
Intended status: Standards Track
Expires: January 3, 2019

C. Do
W. Kolodziejek
J. Chroboczek
IRIF, University of Paris-Diderot
July 2, 2018

Babel Cryptographic Authentication
draft-do-babel-hmac-00

Abstract

This document describes a cryptographic authentication mechanism for the Babel routing protocol that has provisions for replay avoidance. This document updates RFC 6126bis and obsoletes RFC 7298.

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 January 3, 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
1.1. Applicability	3
1.2. Assumptions and security properties	3
1.3. Specification of Requirements	4
2. Conceptual overview of the protocol	4
3. Data Structures	5
3.1. The Interface Table	5
3.2. The Neighbour table	6
4. Protocol Operation	6
4.1. HMAC computation	6
4.2. Packet Transmission	7
4.3. Packet Reception	7
5. Packet Format	10
5.1. HMAC TLV	10
5.2. PC TLV	10
5.3. Challenge Request TLV	11
5.4. Challenge Reply TLV	11
6. Security Considerations	12
7. IANA Considerations	12
8. Acknowledgments	12
9. References	12
9.1. Normative References	12
9.2. Informational References	13
Appendix A. Use of the packet trailer	13
Appendix B. Incremental deployment and key rotation	13
Appendix C. Implicit indices	14
Authors' Addresses	16

1. Introduction

By default, the Babel routing protocol trusts the information contained in every UDP packet it receives on the Babel port. An attacker can redirect traffic to itself or to a different node in the network, causing a variety of potential issues. In particular, an attacker might:

- o spoof a Babel packet, and redirect traffic by announcing a smaller metric, a larger seqno, or a longer prefix;
- o spoof a malformed packet, which could cause an insufficiently robust implementation to crash or interfere with the rest of the network;
- o replay a previously captured Babel packet, which could cause traffic to be redirected or otherwise interfere with the network.

Protecting a Babel network is challenging due to the fact that the Babel protocol uses both unicast and multicast communication. One possible approach, used notably by the Babel over DTLS protocol, is to require a secured version of Babel to use unicast communication for all semantically significant communication, and then use a standard unicast security protocol to protect the Babel traffic. In this document, we take the opposite approach: we define a cryptographic extension to the Babel protocol that is able to protect both unicast and multicast traffic, and thus requires very few changes to the core protocol.

1.1. Applicability

The protocol defined in this document assumes that all interfaces on a given link are equally trusted and share a small set of symmetric keys (usually just one, two during key rotation). The protocol is inapplicable in situations where asymmetric keying is required, where the trust relationship is partial, or where large numbers of trusted keys are provisioned on a single link at the same time.

This protocol supports incremental deployment (where an insecure Babel network is made secure with no service interruption), and it supports graceful key rotation (where the set of keys is changed with no service interruption).

This protocol does not require synchronised clocks, it does not require persistently monotonic clocks, and it does not require any form of persistent storage.

1.2. Assumptions and security properties

The correctness of the protocol relies on the following assumptions:

- o that the HMAC being used is invulnerable to spoofing, i.e. that an attacker is unable to generate a packet with a correct HMAC;
- o that a node never generates the same index or nonce twice over the lifetime of a key.

The first assumption is a property of the HMAC being used, and is therefore out-of-scope for this document. The second assumption can be met either by using a robust random number generator and sufficiently large indices and nonces, by using a reliable hardware clock, or by rekeying whenever a collision becomes likely.

If the assumptions above are met, the protocol described in this document has the following properties:

- o it is invulnerable to spoofing: any packet accepted as authentic is the exact copy of a packet originally sent by an authorised node;
- o locally to a single node, it is invulnerable to replay: if a node has previously accepted a given packet, then it will never again accept a copy of this packet or an earlier packet from the same sender;
- o among different nodes, it is only vulnerable to immediate replay: if a node A has accepted a packet from C as valid, then a node B will only accept a copy of that packet as authentic if B has accepted an older packet from C and B has received no later packet from C.

While this protocol makes serious efforts to mitigate the effects of a denial of service attack, it does not fully protect against such attacks.

1.3. Specification of Requirements

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. Conceptual overview of the protocol

When a node B sends out a Babel packet through an interface that is configured for cryptographic protection, it computes one or more HMACs which it appends to the packet. When a node A receives a packet over an interface that requires cryptographic protection, it independently computes a set of HMACs and compares them to the HMACs appended to the packet; if there is no match, the packet is discarded.

In order to protect against replay B maintains a per-interface 32-bit integer known as the "packet counter" (PC). Whenever B sends a packet through the interface, it embeds the current value of the PC within the region of the packet that is protected by the HMACs and increases the PC by at least one. When A receives the packet, it compares the value of the PC with the one contained in the previous packet received from B, and unless it is strictly greater, the packet is discarded.

By itself, the PC mechanism is not sufficient to protect against replay. Consider a peer A that has no information about a pair B

(e.g., because it has recently rebooted). Suppose that A receives a packet ostensibly from B carrying a given PC; since A has no information about B, it has no way to determine whether the packet is freshly generated or a replay of a previously sent packet.

In this situation, A discards the packet and challenges B to prove that it knows the HMAC key. It sends a "challenge request", a TLV containing a unique nonce, a value that has never been used before and will never be used again. B replies to the challenge request with a "challenge reply", a TLV containing a copy of the nonce chosen by A, in a packet protected by HMAC and containing the new value of B's PC. Since the nonce has never been used before, B's reply proves B's knowledge of the HMAC key and the freshness of the PC.

By itself, this mechanism is safe against replay if B never resets its PC. In practice, however, this is difficult to ensure, as persistent storage is prone to failure, and hardware clocks, even when available, are occasionally reset. Suppose that B resets its PC to an earlier value, and sends a packet with a previously used PC n . A challenges B, B successfully responds to the challenge, and A accepts the PC equal to $n + 1$. At this point, an attacker C may send a replayed packet with PC equal to $n + 2$, which will be accepted by A.

Another mechanism is needed to protect against this attack. In this protocol, every PC is tagged with an "index", an arbitrary string of octets. Whenever B resets its PC, or whenever B doesn't know whether its PC has been reset, it picks an index that it has never used before (either by drawing it randomly or by using a reliable hardware clock) and starts sending PCs with that index. Whenever A detects that B has changed its index, it challenges B again.

With this additional mechanism, this protocol is provably invulnerable to replay attacks (see Section 1.2 above).

3. Data Structures

3.1. The Interface Table

Every Babel node maintains an interface table, as described in [RFC6126bis] Section 3.2.3. This protocol extends the entries in this table with a set of HMAC keys, and a pair (Index, PC), where Index is an arbitrary string of bytes and PC is a 32-bit integer. The Index is initialised to a value that has never been used before (e.g., by choosing a random string of sufficient length).

3.2. The Neighbour table

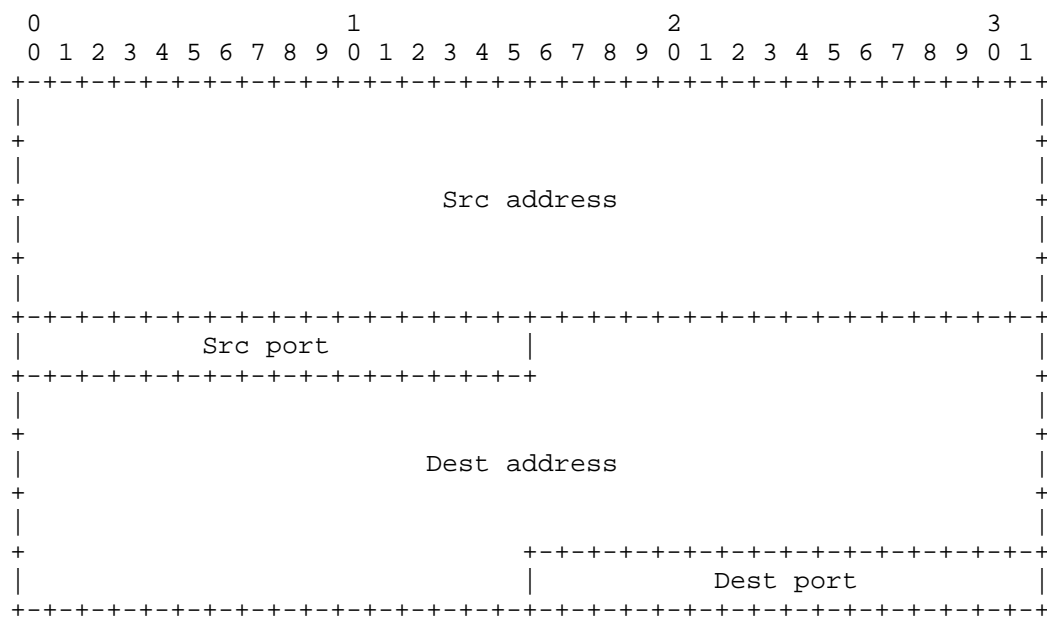
Every Babel node maintains a neighbour table, as described in [RFC6126bis] Section 3.2.4. This protocol extends the entries in this table with a pair (Index, PC), as well as a nonce (an arbitrary string of bytes) and a challenge expiry timer. The Index and PC are initially undefined, and are managed as described in Section 4.3. The Nonce and expiry timer are initially undefined and used as described in Section 4.3.1.1.

4. Protocol Operation

4.1. HMAC computation

A Babel node computes an HMAC as follows.

First, the node builds a pseudo-header that will participate in HMAC computation but will not be sent. The pseudo-header has the following format:



Fields :

Src address The source IP address of the packet.

Src port	The source UDP port number of the packet.
----------	---

Dest address The destination IP address of the packet.

Src port The destination UDP port number of the packet.

The node takes the concatenation of the pseudo-header and the packet including the packet header but excluding the packet trailer (from octet 0 inclusive up to Body Length + 4 exclusive) and computes an HMAC as defined in Section 2 of [RFC2104] with one of the implemented hash algorithms. Every implementation MUST implement HMAC-SHA256 [RFC6234], and MAY implement other HMAC algorithms.

4.2. Packet Transmission

A Babel node may delay actually sending TLVs by a small amount, in order to aggregate multiple TLVs in a single packet up to the interface MTU (Section 4 of [RFC6126bis]). For an interface on which HMAC protection is configured, the TLV aggregation logic MUST take into account the overhead due to PC TLVs (one in each packet) and HMAC TLVs (one per configured key).

Before sending a packet, the following actions are performed:

- o a PC TLV containing the Packet Counter and Index associated with the outgoing interface is appended to the packet body; the packet counter is incremented by a strictly positive amount (typically just 1); if the packet counter overflows, a new index is generated;
- o for each key configured on the interface, an HMAC is computed as specified in Section 4.1 above, and an HMAC TLV is appended to the packet trailer.

4.3. Packet Reception

When a packet is received on an interface that is configured for HMAC protection, the following steps are performed before the packet is passed to normal processing:

- o First, the receiver checks whether the trailer of the received packet carries at least one HMAC TLV; if not, the packet is immediately dropped and processing stops. Then, for each key configured on the receiving interface, the implementation computes the HMAC of the packet. It then compares every generated HMAC against every HMAC included in the packet; if there is at least one match, the packet passes the HMAC test; if there is none, the packet is silently dropped and processing stops at this point. In order to avoid memory exhaustion attacks, an entry in the Neighbour Table MUST NOT be created before the HMAC test has

passed successfully. The HMAC of the packet MUST NOT be computed for each HMAC TLV contained in the packet, but only once for each configured key.

- o The packet body is then parsed a first time. During this "preparse" phase, the packet body is traversed and all TLVs are ignored except PC TLVs, Challenge Requests and Challenge Replies. When a PC TLV is encountered, the enclosed PC and Index are saved for later processing; if multiple PCs are found, only the first one is processed, the remaining ones are silently ignored. If a Challenge Request is encountered, a Challenge Reply is scheduled, as described in Section 4.3.1.2, and if a Challenge Reply is encountered, it is tested for validity as described in Section 4.3.1.3 and a note is made of the result of the test.
- o The preparse phase above has yielded two pieces of data: the PC and Index from the first PC TLV, and a bit indicating whether the packet contains a successful Challenge Reply. If the packet does not contain a PC TLV, the packet is dropped and processing stops at this point. If the packet contains a successful Challenge Reply, then the PC and Index contained in the PC TLV are stored in the Neighbour Table entry corresponding to the sender (which may need to be created at this stage).
- o If there is no entry in the Neighbour Table corresponding to the sender, or if such an entry exists but contains no Index, or if the Index it contains is different from the Index contained in the PC TLV, then a challenge is sent as described in Section 4.3.1.1, processing stops at this stage, and the packet is dropped.
- o At this stage, the Index contained in the PC TLV is equal to the Index in the Neighbour Table entry corresponding to the sender. The receiver compares the received PC with the PC contained in the Neighbour Table; if the received PC smaller or equal than the PC contained in the Neighbour Table, the packet is silently dropped and processing stops (no challenge is sent in this case, since the mismatch might be caused by harmless packet reordering on the link). Otherwise, the PC contained in the Neighbour Table entry is set to the received PC, and the packet is accepted.

After the packet has been accepted, it is processed as normal, except that any PC, Challenge Request and Challenge Reply TLVs that it contains are silently ignored.

4.3.1. Challenge Requests and Replies

During the preparse stage, the receiver might encounter a mismatched Index, to which it will react by scheduling a Challenge Request. It might encounter a Challenge Request TLV, to which it will reply with a Challenge Reply TLV. Finally, it might encounter a Challenge Reply TLV, which it will attempt to match with a previously sent Challenge Request TLV in order to update the Neighbour Table entry corresponding to the sender of the packet.

4.3.1.1. Sending challenges

When it encounters a mismatched Index during the preparse phase, a node picks a nonce that it has never used before, for example by drawing a sufficiently large random string of bytes or by consulting a strictly monotonic hardware clock. It stores the nonce in the entry of the Neighbour Table of the neighbour (the entry might need to be created at this stage), initialises the neighbour's challenge expiry timer to 30 seconds, and sends a Challenge Request TLV to the unicast address corresponding to the neighbour.

A node MAY aggregate a Challenge Request with other TLVs; in other words, if it has already buffered TLVs to be sent to the unicast address of the sender of the neighbour, it MAY send the buffered TLVs in the same packet as the Challenge Request. However, it MUST arrange for the Challenge Request to be sent in a timely manner, as any packets received from that neighbour will be silently ignored until the challenge completes.

Since a challenge may be prompted by a replayed packet, a node MUST impose a rate limitation to the challenges it sends; a limit of one challenge every 300ms for each neighbour is suggested.

4.3.1.2. Replying to challenges

When it encounters a Challenge Request during the preparse phase, a node constructs a Challenge Reply TLV by copying the Nonce from the Challenge Request into the Challenge Reply. It sends the Challenge Reply to the unicast address of the sender of the Challenge Reply.

A node MAY aggregate a Challenge Reply with other TLVs; in other words, if it has already buffered TLVs to be sent to the unicast address of the sender of the Challenge Request, it MAY send the buffered TLVs in the same packet as the Challenge Reply. However, it MUST arrange for the Challenge Reply to be sent in a timely manner (within a few seconds), and SHOULD NOT send any other packets over the same interface before sending the Challenge Reply, as those would be dropped by the challenger.

A challenge sent to a multicast address MUST be silently ignored.

4.3.1.3. Receiving challenge replies

When it encounters a Challenge Reply during the preparse phase, a node consults the Neighbour Table entry corresponding to the neighbour that sent the Challenge Reply. If no challenge is in progress, i.e., if there is no Nonce stored in the Neighbour Table entry or the Challenge timer has expired, the Challenge Reply is silently ignored and the challenge has failed.

Otherwise, the node compares the Nonce contained in the Challenge Reply with the Nonce contained in the Neighbour Table entry. If the two are equal (they have the same length and content), then the challenge has succeeded; otherwise, the challenge has failed.

5. Packet Format

5.1. HMAC TLV

```

      0                               1                               2                               3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      Type      |      Length      |      HMAC...      |
+-----+-----+-----+-----+-----+-----+-----+-----+

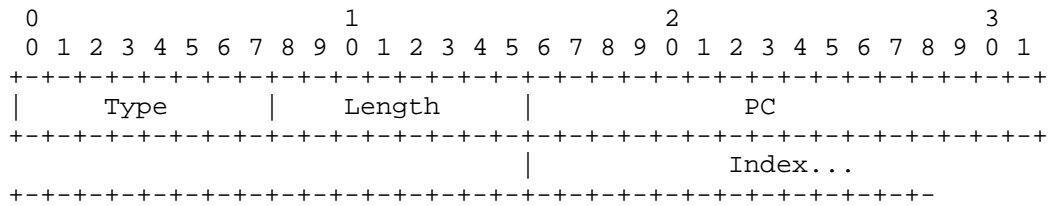
```

Fields :

Type	Set to TBD to indicate an HMAC TLV
Length	The length of the body, exclusive of the Type and Length fields. The length of the body depends on the hash function used.
HMAC	The body contains the HMAC of the whole packet plus the pseudo header.

This TLV is allowed in the packet trailer (see Appendix A), and MUST BE ignored if it is found in the packet body.

5.2. PC TLV



Fields :

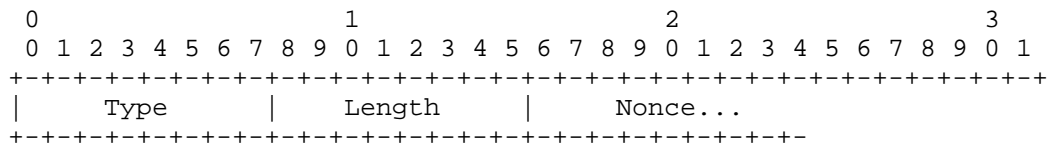
Type Set to TBD to indicate a PC TLV

Length The length of the body, exclusive of the Type and Length fields.

PC The Packet Counter (PC), which is increased with every packet sent over this interface. A new index **MUST** be generated whenever the PC overflows.

Index The sender's Index.

5.3. Challenge Request TLV



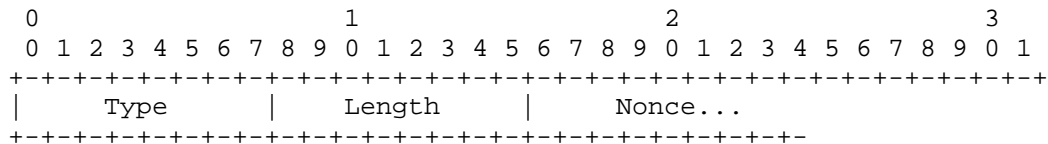
Fields :

Type Set to TBD to indicate a Challenge Request TLV

Length The length of the body, exclusive of the Type and Length fields.

Nonce The nonce uniquely identifying the challenge.

5.4. Challenge Reply TLV



Fields :

Type Set to TBD to indicate a Challenge Reply TLV

Length The length of the body, exclusive of the Type and Length fields. The length of the body is set to the same size as the challenge request TLV length received.

Nonce A copy of the nonce contained in the corresponding challenge request.

6. Security Considerations

7. IANA Considerations

IANA is instructed to allocate the following values in the Babel TLV Numbers registry:

Type	Name	Reference
TBD	HMAC	this document
TBD	PC	this document
TBD	Challenge Request	this document
TBD	Challenge Reply	this document

8. Acknowledgments

The protocol described in this document is based on the original HMAC protocol defined by Denis Ovsienko [RFC7298]. The use of a pseudo-header was suggested by David Schinazi. The use of an index to avoid replay was suggested by Markus Stenberg. The authors are also indebted to Florian Horn and Toke Hoyland-Jorgensen.

9. References

9.1. Normative References

- [RFC2104] Krawczyk, H., Bellare, M., and R. Canetti, "HMAC: Keyed-Hashing for Message Authentication", RFC 2104, DOI 10.17487/RFC2104, February 1997, <<https://www.rfc-editor.org/info/rfc2104>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997.

[RFC6126bis]

Chroboczek, J. and D. Schinazi, "The Babel Routing Protocol", draft-ietf-babel-rfc6126bis-05 (work in progress), May 2018.

[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>>.

[RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017.

9.2. Informational References

[RFC7298] Ovsienko, D., "Babel Hashed Message Authentication Code (HMAC) Cryptographic Authentication", RFC 7298, DOI 10.17487/RFC7298, July 2014, <<https://www.rfc-editor.org/info/rfc7298>>.

Appendix A. Use of the packet trailer

The protocol described in this document uses the packet trailer for storing HMAC TLVs. RFC 6126bis [RFC6126bis] leaves the format of the packet trailer undefined. If the final version of this specification uses the packet trailer, RFC 6126bis will need to be extended with information about the format of the packet trailer.

This document assumes that the packet trailer has the same format as the packet body, i.e., that it consists of a sequence of TLVs. The receiver **MUST** silently ignore any TLV found in the packet trailer unless its definition states that the TLV is allowed in the packet trailer.

Appendix B. Incremental deployment and key rotation

This protocol supports incremental deployment (transitioning from an insecure network to a secured network with no service interruption) and key rotation (transitioning from a set of keys to a different set of keys).

In order to perform incremental deployment, the nodes in the network are first configured in a mode where packets are sent with authentication but not checked on reception. Once all the nodes in the network are configured to send authenticated packets, nodes are reconfigured to reject unauthenticated packets.

In order to perform key rotation, the new key is added to all the nodes; once this is done, both the old and the new key are sent in all packets, and packets are accepted if they are properly signed by either of the keys. At that point, the old key is removed.

In order to support incremental deployment and key rotation, implementations SHOULD support an interface configuration in which they send authenticated packets but accept all packets, and SHOULD allow changing the set of keys associated with an interface without a restart.

Appendix C. Implicit indices

[This appendix describes the "implicit indices" variant of the protocol, which is different and incompatible to the "explicit indices" variant described in the body of this document. This section should either be integrated into the body of the document or removed before publication of this document as an RFC, depending on which protocol variant is finally chosen.]

The protocol described in the body of this document explicitly sends indices as in each packet as part of the PC TLV. Observe that, except when a challenge is required, the index sent on the wire is identical to the index stored in the Neighbour Table, and therefore doesn't need to be sent explicitly except during challenges: it is enough for it to participate in HMAC computation in order to protect against replay. The "implicit indices" variant of the protocol, due to Markus Stenberg and described in this appendix, uses this observation to avoid sending indices explicitly and thus shaves off 2 to 16 octets from almost every packet.

The changes to the protocol are as follows. The pseudo-header includes the Index, and therefore has the following format:

the packet body. In any case, it uses its current index to generate the pseudo-header that will be used to compute the HMAC. (This implies that a packet must be parsed in its entirety before HMAC validation, which requires a robust parser.)

Packet reception is modified as follows. Before checking the HMAC of a packet, the receiver checks whether the packet contains an explicit index. If this is the case, it uses the index contained in the packet in order to generate the pseudo header; if this is not the case, it uses the index contained in its neighbours table. If there is no index available for that neighbour (either because the table doesn't contain an entry for this neighbour, or the entry doesn't contain an index), HMAC validation fails.

The index and PC contained in the neighbours table are only updated after HMAC validation has succeeded.

Since it is now impossible to differentiate between a failed HMAC and an index change, a node must send a challenge whenever HMAC validation fails. This implies that spoofed packets cause a spurious challenge, but that doesn't change the security properties of the protocol much, given that in any case replayed packets can be used to cause a spurious challenge.

Authors' Addresses

Clara Do
IRIF, University of Paris-Diderot
75205 Paris Cedex 13
France

Email: clarado_perso@yahoo.fr

Weronika Kolodziejak
IRIF, University of Paris-Diderot
75205 Paris Cedex 13
France

Email: weronika.kolodziejak@gmail.com

Juliusz Chroboczek
IRIF, University of Paris-Diderot
Case 7014
75205 Paris Cedex 13
France

Email: jch@irif.fr

Babel routing protocol
Internet-Draft
Intended status: Informational
Expires: 12 September 2021

B.H. Stark
AT&T
M.J. Jethanandani
VMware
11 March 2021

Babel Information Model
draft-ietf-babel-information-model-14

Abstract

This Babel Information Model provides structured data elements for a Babel implementation reporting its current state and may allow limited configuration of some such data elements. This information model can be used as a basis for creating data models under various data modeling regimes. This information model only includes parameters and parameter values useful for managing Babel over IPv6.

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 12 September 2021.

Copyright Notice

Copyright (c) 2021 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
1.1. Requirements Language	3
1.2. Notation	3
2. Overview	4
3. The Information Model	7
3.1. Definition of babel-information-obj	7
3.2. Definition of babel-constants-obj	9
3.3. Definition of babel-interface-obj	9
3.4. Definition of babel-if-stats-obj	12
3.5. Definition of babel-neighbor-obj	13
3.6. Definition of babel-route-obj	14
3.7. Definition of babel-mac-key-set-obj	16
3.8. Definition of babel-mac-key-obj	16
3.9. Definition of babel-dtls-cert-set-obj	18
3.10. Definition of babel-dtls-cert-obj	18
4. Extending the Information Model	19
5. Security Considerations	19
6. IANA Considerations	20
7. Acknowledgements	20
8. References	21
8.1. Normative References	21
8.2. Informative References	22
Authors' Addresses	22

1. Introduction

Babel is a loop-avoiding distance-vector routing protocol defined in [RFC8966]. [RFC8967] defines a security mechanism that allows Babel packets to be cryptographically authenticated, and [RFC8968] defines a security mechanism that allows Babel packets to be both authenticated and encrypted. This document describes an information model for Babel (including implementations using one or both of these security mechanisms) that can be used to create management protocol data models (such as a NETCONF [RFC6241] YANG [RFC7950] data model).

Due to the simplicity of the Babel protocol, most of the information model is focused on reporting Babel protocol operational state, and very little of that is considered mandatory to implement for an implementation claiming compliance with this information model. Some parameters may be configurable. However, it is up to the Babel implementation whether to allow any of these to be configured within its implementation. Where the implementation does not allow configuration of these parameters, it MAY still choose to expose them as read-only.

The Information Model is presented using a hierarchical structure. This does not preclude a data model based on this Information Model from using a referential or other structure.

This information model only includes parameters and parameter values useful for managing Babel over IPv6. This model has no parameters or values specific to operating Babel over IPv4, even though [RFC8966] does define a multicast group for sending and listening to multicast announcements on IPv4. There is less likelihood of breakage due to inconsistent configuration and increased implementation simplicity if Babel is operated always and only over IPv6. Running Babel over IPv6 requires IPv6 at the link layer and does not need advertised prefixes, router advertisements or DHCPv6 to be present in the network. Link-local IPv6 is widely supported among devices where Babel is expected to be used. Note that Babel over IPv6 can be used for configuration of both IPv4 and IPv6 routes.

1.1. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP014 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

1.2. Notation

This document uses a programming language-like notation to define the properties of the objects of the information model. An optional property is enclosed by square brackets, [], and a list property is indicated by two numbers in angle brackets, <m..n>, where m indicates the minimal number of list elements, and n indicates the maximum number of list elements. The symbol * for n means there are no defined limits on the number of list elements. Each parameter and object includes an indication of "ro" or "rw". "ro" means the parameter or object is read-only. "rw" means it is read-write. For an object, read-write means instances of the object can be created or deleted. If an implementation is allowed to choose to implement a

"rw" parameter as read-only, this is noted in the parameter description.

The object definitions use base types that are defined as follows:

binary	A binary string (sequence of octets).
boolean	A type representing a Boolean (true or false) value.
datetime	A type representing a date and time using the Gregorian calendar. The datetime format MUST conform to RFC 3339 [RFC3339] Section 5.6.
ip-address	A type representing an IP address. This type supports both IPv4 and IPv6 addresses.
operation	A type representing a remote procedure call or other action that can be used to manipulate data elements or system behaviors.
reference	A type representing a reference to another information or data model element or to some other device resource.
string	A type representing a human-readable string consisting of a (possibly restricted) subset of Unicode and ISO/IEC 10646 [ISO.10646] characters.
uint	A type representing an unsigned integer number. This information model does not define a precision.

2. Overview

The Information Model is hierarchically structured as follows:

```
+-- babel-information
  +-- babel-implementation-version
  +-- babel-enable
  +-- router-id
  +-- self-seqno
  +-- babel-metric-comp-algorithms
  +-- babel-security-supported
  +-- babel-mac-algorithms
  +-- babel-dtls-cert-types
  +-- babel-stats-enable
  +-- babel-stats-reset
  +-- babel-constants
    | +-- babel-udp-port
    | +-- babel-mcast-group
```

```
+--- babel-interfaces
|   +--- babel-interface-reference
|   +--- babel-interface-enable
|   +--- babel-interface-metric-algorithm
|   +--- babel-interface-split-horizon
|   +--- babel-mcast-hello-seqno
|   +--- babel-mcast-hello-interval
|   +--- babel-update-interval
|   +--- babel-mac-enable
|   +--- babel-if-mac-key-sets
|   +--- babel-mac-verify
|   +--- babel-dtls-enable
|   +--- babel-if-dtls-cert-sets
|   +--- babel-dtls-cached-info
|   +--- babel-dtls-cert-prefer
|   +--- babel-packet-log-enable
|   +--- babel-packet-log
|   +--- babel-if-stats
|       +--- babel-sent-mcast-hello
|       +--- babel-sent-mcast-update
|       +--- babel-sent-ucast-hello
|       +--- babel-sent-ucast-update
|       +--- babel-sent-IHU
|       +--- babel-received-packets
|   +--- babel-neighbors
|       +--- babel-neighbor-address
|       +--- babel-hello-mcast-history
|       +--- babel-hello-ucast-history
|       +--- babel-txcost
|       +--- babel-exp-mcast-hello-seqno
|       +--- babel-exp-ucast-hello-seqno
|       +--- babel-ucast-hello-seqno
|       +--- babel-ucast-hello-interval
|       +--- babel-rxcost
|       +--- babel-cost
+--- babel-routes
|   +--- babel-route-prefix
|   +--- babel-route-prefix-length
|   +--- babel-route-router-id
|   +--- babel-route-neighbor
|   +--- babel-route-received-metric
|   +--- babel-route-calculated-metric
|   +--- babel-route-seqno
|   +--- babel-route-next-hop
|   +--- babel-route-feasible
|   +--- babel-route-selected
+--- babel-mac-key-sets
|   +--- babel-mac-default-apply
```

```
|  +-- babel-mac-keys
|  |  +-- babel-mac-key-name
|  |  +-- babel-mac-key-use-send
|  |  +-- babel-mac-key-use-verify
|  |  +-- babel-mac-key-value
|  |  +-- babel-mac-key-algorithm
|  |  +-- babel-mac-key-test
+-- babel-dtls-cert-sets
  +-- babel-dtls-default-apply
  +-- babel-dtls-certs
    +-- babel-cert-name
    +-- babel-cert-value
    +-- babel-cert-type
    +-- babel-cert-private-key
```

Most parameters are read-only. Following is a descriptive list of the parameters that are not required to be read-only:

- * enable/disable Babel
- * create/delete Babel MAC Key sets
- * create/delete Babel Certificate sets
- * enable/disable statistics collection
- * Constant: UDP port
- * Constant: IPv6 multicast group
- * Interface: enable/disable Babel on this interface
- * Interface: Metric algorithm
- * Interface: Split horizon
- * Interface: sets of MAC keys
- * Interface: verify received MAC packets
- * Interface: set of certificates for use with DTLS
- * Interface: use cached info extensions
- * Interface: preferred order of certificate types
- * Interface: enable/disable packet log

- * MAC-keys: create/delete entries
- * MAC-keys: key used for sent packets
- * MAC-keys: key used to verify packets
- * DTLS-certs: create/delete entries

The following parameters are required to return no value when read:

- * MAC key values
- * DTLS private keys

Note that this overview is intended simply to be informative and is not normative. If there is any discrepancy between this overview and the detailed information model definitions in subsequent sections, the error is in this overview.

3. The Information Model

3.1. Definition of babel-information-obj

```
object {
    string                ro babel-implementation-version;
    boolean               rw babel-enable;
    binary                ro babel-self-router-id;
    [uint                 ro babel-self-seqno;]
    string                ro babel-metric-comp-algorithms<1..*>;
    string                ro babel-security-supported<0..*>;
    [string               ro babel-mac-algorithms<1..*>;]
    [string               ro babel-dtls-cert-types<1..*>;]
    [boolean              rw babel-stats-enable;]
    [operation            babel-stats-reset;]
    babel-constants-obj   ro babel-constants;
    babel-interface-obj   ro babel-interfaces<0..*>;
    babel-route-obj       ro babel-routes<0..*>;
    [babel-mac-key-set-obj rw babel-mac-key-sets<0..*>;]
    [babel-dtls-cert-set-obj rw babel-dtls-cert-sets<0..*>;]
} babel-information-obj;
```

babel-implementation-version: The name and version of this implementation of the Babel protocol.

babel-enable: When written, it configures whether the protocol should be enabled (true) or disabled (false). A read from the running or intended datastore indicates the configured administrative value of whether the protocol is enabled (true) or

not (false). A read from the operational datastore indicates whether the protocol is actually running (true) or not (i.e., it indicates the operational state of the protocol). A data model that does not replicate parameters for running and operational datastores can implement this as two separate parameters. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-self-router-id: The router-id used by this instance of the Babel protocol to identify itself. [RFC8966] describes this as an arbitrary string of 8 octets.

babel-self-seqno: The current sequence number included in route updates for routes originated by this node. This is a 16-bit unsigned integer.

babel-metric-comp-algorithms: List of supported cost computation algorithms. Possible values include "2-out-of-3", and "ETX". "2-out-of-3" is described in [RFC8966], section A.2.1. "ETX" is described in [RFC8966], section A.2.2.

babel-security-supported: List of supported security mechanisms. Possible values include "MAC" to indicate support of [RFC8967] and "DTLS" to indicate support of [RFC8968].

babel-mac-algorithms: List of supported MAC computation algorithms. Possible values include "HMAC-SHA256", "BLAKE2s-128" to indicate support for algorithms indicated in [RFC8967].

babel-dtls-cert-types: List of supported certificate types. Possible values include "X.509" and "RawPublicKey" to indicate support for types indicated in [RFC8968].

babel-stats-enable: Indicates whether statistics collection is enabled (true) or disabled (false) on all interfaces. When enabled, existing statistics values are not cleared and will be incremented as new packets are counted.

babel-stats-reset: An operation that resets all babel-if-stats parameters to zero. This operation has no input or output parameters.

babel-constants: A babel-constants-obj object.

babel-interfaces: A set of babel-interface-obj objects.

babel-routes: A set of babel-route-obj objects. Contains the routes known to this node.

babel-mac-key-sets: A set of babel-mac-key-set-obj objects. If this object is implemented, it provides access to parameters related to the MAC security mechanism. An implementation MAY choose to expose this object as read-only ("ro").

babel-dtls-cert-sets: A set of babel-dtls-cert-set-obj objects. If this object is implemented, it provides access to parameters related to the DTLS security mechanism. An implementation MAY choose to expose this object as read-only ("ro").

3.2. Definition of babel-constants-obj

```
object {  
    uint          rw babel-udp-port;  
    [ip-address    rw babel-mcast-group;]  
} babel-constants-obj;
```

babel-udp-port: UDP port for sending and listening for Babel packets. Default is 6696. An implementation MAY choose to expose this parameter as read-only ("ro"). This is a 16-bit unsigned integer.

babel-mcast-group: Multicast group for sending and listening to multicast announcements on IPv6. Default is ff02::1:6. An implementation MAY choose to expose this parameter as read-only ("ro").

3.3. Definition of babel-interface-obj

```

object {
  reference          ro babel-interface-reference;
  [boolean           rw babel-interface-enable;]
  string            rw babel-interface-metric-algorithm;
  [boolean           rw babel-interface-split-horizon;]
  [uint             ro babel-mcast-hello-seqno;]
  [uint             ro babel-mcast-hello-interval;]
  [uint             ro babel-update-interval;]
  [boolean           rw babel-mac-enable;]
  [reference         rw babel-if-mac-key-sets<0..*>;]
  [boolean           rw babel-mac-verify;]
  [boolean           rw babel-dtls-enable;]
  [reference         rw babel-if-dtls-cert-sets<0..*>;]
  [boolean           rw babel-dtls-cached-info;]
  [string           rw babel-dtls-cert-prefer<0..*>;]
  [boolean           rw babel-packet-log-enable;]
  [reference         ro babel-packet-log;]
  [babel-if-stats-obj ro babel-if-stats;]
  babel-neighbor-obj ro babel-neighbors<0..*>;
} babel-interface-obj;

```

babel-interface-reference: Reference to an interface object that can be used to send and receive IPv6 packets, as defined by the data model (e.g., YANG [RFC7950], BBF [TR-181]). Referencing syntax will be specific to the data model. If there is no set of interface objects available, this should be a string that indicates the interface name used by the underlying operating system.

babel-interface-enable: When written, it configures whether the protocol should be enabled (true) or disabled (false) on this interface. A read from the running or intended datastore indicates the configured administrative value of whether the protocol is enabled (true) or not (false). A read from the operational datastore indicates whether the protocol is actually running (true) or not (i.e., it indicates the operational state of the protocol). A data model that does not replicate parameters for running and operational datastores can implement this as two separate parameters. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-interface-metric-algorithm: Indicates the metric computation algorithm used on this interface. The value MUST be one of those listed in the babel-information-obj babel-metric-comp-algorithms parameter. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-interface-split-horizon: Indicates whether or not the split

horizon optimization is used when calculating metrics on this interface. A value of true indicates split horizon optimization is used. Split horizon optimization is described in [RFC8966], section 3.7.4. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-mcast-hello-seqno: The current sequence number in use for multicast Hellos sent on this interface. This is a 16-bit unsigned integer.

babel-mcast-hello-interval: The current interval in use for multicast Hellos sent on this interface. Units are centiseconds. This is a 16-bit unsigned integer.

babel-update-interval: The current interval in use for all updates (multicast and unicast) sent on this interface. Units are centiseconds. This is a 16-bit unsigned integer.

babel-mac-enable: Indicates whether the MAC security mechanism is enabled (true) or disabled (false). An implementation MAY choose to expose this parameter as read-only ("ro").

babel-if-mac-keys-sets: List of references to the babel-mac entries that apply to this interface. When an interface instance is created, all babel-mac-key-sets instances with babel-mac-default-apply "true" will be included in this list. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-mac-verify A Boolean flag indicating whether MACs in incoming Babel packets are required to be present and are verified. If this parameter is "true", incoming packets are required to have a valid MAC. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-dtls-enable: Indicates whether the DTLS security mechanism is enabled (true) or disabled (false). An implementation MAY choose to expose this parameter as read-only ("ro").

babel-if-dtls-cert-sets: List of references to the babel-dtls-cert-sets entries that apply to this interface. When an interface instance is created, all babel-dtls-cert-sets instances with babel-dtls-default-apply "true" will be included in this list. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-dtls-cached-info: Indicates whether the cached_info extension

(see [RFC8968] Appendix A) is included in ClientHello and ServerHello packets. The extension is included if the value is "true". An implementation MAY choose to expose this parameter as read-only ("ro").

babel-dtls-cert-prefer: List of supported certificate types, in order of preference. The values MUST be among those listed in the babel-dtls-cert-types parameter. This list is used to populate the server_certificate_type extension (see [RFC8968] Appendix A) in a Client Hello. Values that are present in at least one instance in the babel-dtls-certs object of a referenced babel-dtls instance and that have a non-empty babel-cert-private-key will be used to populate the client_certificate_type extension in a Client Hello.

babel-packet-log-enable: Indicates whether packet logging is enabled (true) or disabled (false) on this interface.

babel-packet-log: A reference or url link to a file that contains a timestamped log of packets received and sent on babel-udp-port on this interface. The [libpcap] file format with .pcap file extension SHOULD be supported for packet log files. Logging is enabled / disabled by babel-packet-log-enable. Implementations will need to carefully manage and limit memory used by packet logs.

babel-if-stats: Statistics collection object for this interface.

babel-neighbors: A set of babel-neighbor-obj objects.

3.4. Definition of babel-if-stats-obj

```
object {  
    uint    ro babel-sent-mcast-hello;  
    uint    ro babel-sent-mcast-update;  
    uint    ro babel-sent-ucast-hello;  
    uint    ro babel-sent-ucast-update;  
    uint    ro babel-sent-IHU;  
    uint    ro babel-received-packets;  
} babel-if-stats-obj;
```

babel-sent-mcast-hello: A count of the number of multicast Hello packets sent on this interface.

babel-sent-mcast-update: A count of the number of multicast update packets sent on this interface.

babel-sent-ucast-hello: A count of the number of unicast Hello

packets sent on this interface.

babel-sent-ucast-update: A count of the number of unicast update packets sent on this interface.

babel-sent-IHU: A count of the number of IHU packets sent on this interface.

babel-received-packets: A count of the number of Babel packets received on this interface.

3.5. Definition of babel-neighbor-obj

```
object {  
    ip-address      ro babel-neighbor-address;  
    [binary         ro babel-hello-mcast-history;]  
    [binary         ro babel-hello-ucast-history;]  
    uint            ro babel-txcost;  
    uint            ro babel-exp-mcast-hello-seqno;  
    uint            ro babel-exp-ucast-hello-seqno;  
    [uint           ro babel-ucast-hello-seqno;]  
    [uint           ro babel-ucast-hello-interval;]  
    [uint           ro babel-rxcost;]  
    [uint           ro babel-cost;]  
} babel-neighbor-obj;
```

babel-neighbor-address: IPv4 or IPv6 address the neighbor sends packets from.

babel-hello-mcast-history: The multicast Hello history of whether or not the multicast Hello packets prior to babel-exp-mcast-hello-seqno were received. A binary sequence where the most recently received Hello is expressed as a "1" placed in the left-most bit, with prior bits shifted right (and "0" bits placed between prior Hello bits and most recent Hello for any not-received Hellos). This value should be displayed using hex digits ([0-9a-fA-F]). See [RFC8966], section A.1.

babel-hello-ucast-history: The unicast Hello history of whether or not the unicast Hello packets prior to babel-exp-ucast-hello-seqno were received. A binary sequence where the most recently received Hello is expressed as a "1" placed in the left-most bit, with prior bits shifted right (and "0" bits placed between prior Hello bits and most recent Hello for any not-received Hellos). This value should be displayed using hex digits ([0-9a-fA-F]). See [RFC8966], section A.1.

babel-txcost: Transmission cost value from the last IHU packet

received from this neighbor, or maximum value to indicate the IHU hold timer for this neighbor has expired. See [RFC8966], section 3.4.2. This is a 16-bit unsigned integer.

babel-exp-mcast-hello-seqno: Expected multicast Hello sequence number of next Hello to be received from this neighbor. If multicast Hello packets are not expected, or processing of multicast packets is not enabled, this MUST be NULL. This is a 16-bit unsigned integer; if the data model uses zero (0) to represent NULL values for unsigned integers, the data model MAY use a different data type that allows differentiation between zero (0) and NULL.

babel-exp-ucast-hello-seqno: Expected unicast Hello sequence number of next Hello to be received from this neighbor. If unicast Hello packets are not expected, or processing of unicast packets is not enabled, this MUST be NULL. This is a 16-bit unsigned integer; if the data model uses zero (0) to represent NULL values for unsigned integers, the data model MAY use a different data type that allows differentiation between zero (0) and NULL.

babel-ucast-hello-seqno: The current sequence number in use for unicast Hellos sent to this neighbor. If unicast Hellos are not being sent, this MUST be NULL. This is a 16-bit unsigned integer; if the data model uses zero (0) to represent NULL values for unsigned integers, the data model MAY use a different data type that allows differentiation between zero (0) and NULL.

babel-ucast-hello-interval: The current interval in use for unicast Hellos sent to this neighbor. Units are centiseconds. This is a 16-bit unsigned integer.

babel-rxcost: Reception cost calculated for this neighbor. This value is usually derived from the Hello history, which may be combined with other data, such as statistics maintained by the link layer. The rxcost is sent to a neighbor in each IHU. See [RFC8966], section 3.4.3. This is a 16-bit unsigned integer.

babel-cost: The link cost, as computed from the values maintained in the neighbor table: the statistics kept in the neighbor table about the reception of Hellos, and the txcost computed from received IHU packets. This is a 16-bit unsigned integer.

3.6. Definition of babel-route-obj

```
object {  
  ip-address    ro babel-route-prefix;  
  uint          ro babel-route-prefix-length;  
  binary        ro babel-route-router-id;  
  reference      ro babel-route-neighbor;  
  uint          ro babel-route-received-metric;  
  uint          ro babel-route-calculated-metric;  
  uint          ro babel-route-seqno;  
  ip-address    ro babel-route-next-hop;  
  boolean       ro babel-route-feasible;  
  boolean       ro babel-route-selected;  
} babel-route-obj;
```

babel-route-prefix: Prefix (expressed in IP address format) for which this route is advertised.

babel-route-prefix-length: Length of the prefix for which this route is advertised.

babel-route-router-id: The router-id of the router that originated this route.

babel-route-neighbor: Reference to the babel-neighbors entry for the neighbor that advertised this route.

babel-route-received-metric: The metric with which this route was advertised by the neighbor, or maximum value to indicate the route was recently retracted and is temporarily unreachable (see Section 3.5.5 of [RFC8966]). This metric will be NULL if the route was not received from a neighbor but was generated through other means. At least one of babel-route-calculated-metric and babel-route-received-metric MUST be non-NULL. Having both be non-NULL is expected for a route that is received and subsequently advertised. This is a 16-bit unsigned integer; if the data model uses zero (0) to represent NULL values for unsigned integers, the data model MAY use a different data type that allows differentiation between zero (0) and NULL.

babel-route-calculated-metric: A calculated metric for this route.

How the metric is calculated is implementation-specific. Maximum value indicates the route was recently retracted and is temporarily unreachable (see Section 3.5.5 of [RFC8966]). At least one of `babel-route-calculated-metric` and `babel-route-received-metric` MUST be non-NULL. Having both be non-NULL is expected for a route that is received and subsequently advertised. This is a 16-bit unsigned integer; if the data model uses zero (0) to represent NULL values for unsigned integers, the data model MAY use a different data type that allows differentiation between zero (0) and NULL.

`babel-route-seqno`: The sequence number with which this route was advertised. This is a 16-bit unsigned integer.

`babel-route-next-hop`: The next-hop address of this route. This will be empty if this route has no next-hop address.

`babel-route-feasible`: A Boolean flag indicating whether this route is feasible, as defined in Section 3.5.1 of [RFC8966]).

`babel-route-selected`: A Boolean flag indicating whether this route is selected (i.e., whether it is currently being used for forwarding and is being advertised).

3.7. Definition of `babel-mac-key-set-obj`

```
object {  
    boolean          rw babel-mac-default-apply;  
    babel-mac-key-obj rw babel-mac-keys<0..*>;  
} babel-mac-key-set-obj;
```

`babel-mac-default-apply`: A Boolean flag indicating whether this object instance is applied to all new `babel-interface` instances, by default. If "true", this instance is applied to new `babel-interfaces` instances at the time they are created, by including it in the `babel-if-mac-key-sets` list. If "false", this instance is not applied to new `babel-interfaces` instances when they are created. An implementation MAY choose to expose this parameter as read-only ("ro").

`babel-mac-keys`: A set of `babel-mac-key-obj` objects.

3.8. Definition of `babel-mac-key-obj`

```
object {  
    string      rw babel-mac-key-name;  
    boolean     rw babel-mac-key-use-send;  
    boolean     rw babel-mac-key-use-verify;  
    binary      -- babel-mac-key-value;  
    string      rw babel-mac-key-algorithm;  
    [operation   babel-mac-key-test;]  
} babel-mac-key-obj;
```

babel-mac-key-name: A unique name for this MAC key that can be used to identify the key in this object instance, since the key value is not allowed to be read. This value **MUST NOT** be empty and can only be provided when this instance is created (i.e., it is not subsequently writable). The value **MAY** be auto-generated if not explicitly supplied when the instance is created.

babel-mac-key-use-send: Indicates whether this key value is used to compute a MAC and include that MAC in the sent Babel packet. A MAC for sent packets is computed using this key if the value is "true". If the value is "false", this key is not used to compute a MAC to include in sent Babel packets. An implementation **MAY** choose to expose this parameter as read-only ("ro").

babel-mac-key-use-verify: Indicates whether this key value is used to verify incoming Babel packets. This key is used to verify incoming packets if the value is "true". If the value is "false", no MAC is computed from this key for comparing with the MAC in an incoming packet. An implementation **MAY** choose to expose this parameter as read-only ("ro").

babel-mac-key-value: The value of the MAC key. An implementation **MUST NOT** allow this parameter to be read. This can be done by always providing an empty string when read, or through permissions, or other means. This value **MUST** be provided when this instance is created, and is not subsequently writable. This value is of a length suitable for the associated babel-mac-key-algorithm. If the algorithm is based on the HMAC construction [RFC2104], the length **MUST** be between 0 and an upper limit that is at least the size of the output length (where "HMAC-SHA256" output length is 32 octets as described in [RFC4868]). Longer lengths **MAY** be supported but are not necessary if the management system has the ability to generate a suitably random value (e.g., by randomly generating a value or by using a key derivation technique as recommended in [RFC8967] Security Considerations). If the algorithm is "BLAKE2s-128", the length **MUST** be between 0 and 32 bytes inclusive as specified by [RFC7693].

babel-mac-key-algorithm The name of the MAC algorithm used with this

key. The value MUST be the same as one of the enumerations listed in the babel-mac-algorithms parameter. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-mac-key-test: An operation that allows the MAC key and MAC algorithm to be tested to see if they produce an expected outcome. Input to this operation are a binary string and a calculated MAC (also in the format of a binary string) for the binary string. The implementation is expected to create a MAC over the binary string using the babel-mac-key-value and the babel-mac-key-algorithm. The output of this operation is a Boolean indication that the calculated MAC matched the input MAC (true) or the MACs did not match (false).

3.9. Definition of babel-dtls-cert-set-obj

```
object {  
    boolean          rw babel-dtls-default-apply;  
    babel-dtls-cert-obj  rw babel-dtls-certs<0..*>;  
} babel-dtls-cert-set-obj;
```

babel-dtls-default-apply: A Boolean flag indicating whether this object instance is applied to all new babel-interface instances, by default. If "true", this instance is applied to new babel-interfaces instances at the time they are created, by including it in the babel-interface-dtls-certs list. If "false", this instance is not applied to new babel-interfaces instances when they are created. An implementation MAY choose to expose this parameter as read-only ("ro").

babel-dtls-certs: A set of babel-dtls-cert-obj objects. This contains both certificates for this implementation to present for authentication, and to accept from others. Certificates with a non-empty babel-cert-private-key can be presented by this implementation for authentication.

3.10. Definition of babel-dtls-cert-obj

```
object {  
    string          rw babel-cert-name;  
    string          rw babel-cert-value;  
    string          rw babel-cert-type;  
    binary          -- babel-cert-private-key;  
} babel-dtls-cert-obj;
```

babel-cert-name: A unique name for this certificate that can be used

to identify the certificate in this object instance, since the value is too long to be useful for identification. This value MUST NOT be empty and can only be provided when this instance is created (i.e., it is not subsequently writable). The value MAY be auto-generated if not explicitly supplied when the instance is created.

babel-cert-value: The certificate in PEM format [RFC7468]. This value MUST be provided when this instance is created, and is not subsequently writable.

babel-cert-type: The name of the certificate type of this object instance. The value MUST be the same as one of the enumerations listed in the `babel-dtls-cert-types` parameter. This value can only be provided when this instance is created, and is not subsequently writable.

babel-cert-private-key: The value of the private key. If this is non-empty, this certificate can be used by this implementation to provide a certificate during DTLS handshaking. An implementation MUST NOT allow this parameter to be read. This can be done by always providing an empty string when read, or through permissions, or other means. This value can only be provided when this instance is created, and is not subsequently writable.

4. Extending the Information Model

Implementations MAY extend this information model with other parameters or objects. For example, an implementation MAY choose to expose Babel route filtering rules by adding a route filtering object with parameters appropriate to how route filtering is done in that implementation. The precise means used to extend the information model would be specific to the data model the implementation uses to expose this information.

5. Security Considerations

This document defines a set of information model objects and parameters that may be exposed to be visible from other devices, and some of which may be configured. Securing access to and ensuring the integrity of this data is in scope of and the responsibility of any data model derived from this information model. Specifically, any YANG [RFC7950] data model is expected to define security exposure of the various parameters, and a [TR-181] data model will be secured by the mechanisms defined for the management protocol used to transport it.

Misconfiguration (whether unintentional or malicious) can prevent reachability or cause poor network performance (increased latency, jitter, etc.). Misconfiguration of security credentials can cause a denial of service condition for the Babel routing protocol. The information in this model discloses network topology, which can be used to mount subsequent attacks on traffic traversing the network.

This information model defines objects that can allow credentials (for this device, for trusted devices, and for trusted certificate authorities) to be added and deleted. Public keys may be exposed through this model. This model requires that private keys and MAC keys never be exposed. Certificates used by [RFC8968] implementations use separate parameters to model the public parts (including the public key) and the private key.

MAC keys are allowed to be as short as zero-length. This is useful for testing. Network operators are RECOMMENDED to follow current best practices for key length and generation of keys related to the MAC algorithm associated with the key. Short (and zero-length) keys are highly susceptible to brute force attacks and therefore SHOULD NOT be used. See the Security Considerations section of [RFC8967] for additional considerations related to MAC keys. The fifth paragraph of [RFC8967] Security Considerations makes some specific key value recommendations that should be noted. It says that if it is necessary to derive keys from a human-readable passphrase, "only the derived keys should be communicated to the routers" and "the original passphrase itself should be kept on the host used to perform the key generation" (which would be the management system in the case of a remote management protocol). It also recommends that keys "should have a length of 32 octets (both for HMAC-SHA256 and BLAKE2s), and be chosen randomly".

This information model uses key sets and certification sets to provide a means of grouping keys and certificates. This makes it easy to use a different set per interface, the same set for one or more interfaces, have a default set in case a new interface is instantiated and to change keys and certificates as needed.

6. IANA Considerations

This document has no IANA actions.

7. Acknowledgements

Juliusz Chroboczek, Toke Hoeiland-Joergensen, David Schinazi, Antonin Decimo, Acee Lindem, and Carsten Bormann have been very helpful in refining this information model.

The language in the Notation section was mostly taken from [RFC8193].

8. References

8.1. Normative References

- [ISO.10646] International Organization for Standardization, "Information Technology – Universal Multiple-Octet Coded Character Set (UCS)", ISO Standard 10646:2014, 2014.
- [libpcap] Wireshark, "Libpcap File Format", 2015, <<https://wiki.wireshark.org/Development/LibpcapFileFormat>>.
- [RFC2104] Krawczyk, H., Bellare, M., and R. Canetti, "HMAC: Keyed-Hashing for Message Authentication", RFC 2104, DOI 10.17487/RFC2104, February 1997, <<https://www.rfc-editor.org/info/rfc2104>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC3339] Klyne, G. and C. Newman, "Date and Time on the Internet: Timestamps", RFC 3339, DOI 10.17487/RFC3339, July 2002, <<https://www.rfc-editor.org/info/rfc3339>>.
- [RFC4868] Kelly, S. and S. Frankel, "Using HMAC-SHA-256, HMAC-SHA-384, and HMAC-SHA-512 with IPsec", RFC 4868, DOI 10.17487/RFC4868, May 2007, <<https://www.rfc-editor.org/info/rfc4868>>.
- [RFC7468] Josefsson, S. and S. Leonard, "Textual Encodings of PKIX, PKCS, and CMS Structures", RFC 7468, DOI 10.17487/RFC7468, April 2015, <<https://www.rfc-editor.org/info/rfc7468>>.
- [RFC7693] Saarinen, M-J., Ed. and J-P. Aumasson, "The BLAKE2 Cryptographic Hash and Message Authentication Code (MAC)", RFC 7693, DOI 10.17487/RFC7693, November 2015, <<https://www.rfc-editor.org/info/rfc7693>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

- [RFC8966] Chroboczek, J. and D. Schinazi, "The Babel Routing Protocol", RFC 8966, DOI 10.17487/RFC8966, January 2021, <<https://www.rfc-editor.org/info/rfc8966>>.
- [RFC8967] Dô, C., Kolodziejak, W., and J. Chroboczek, "MAC Authentication for the Babel Routing Protocol", RFC 8967, DOI 10.17487/RFC8967, January 2021, <<https://www.rfc-editor.org/info/rfc8967>>.
- [RFC8968] Décimo, A., Schinazi, D., and J. Chroboczek, "Babel Routing Protocol over Datagram Transport Layer Security", RFC 8968, DOI 10.17487/RFC8968, January 2021, <<https://www.rfc-editor.org/info/rfc8968>>.

8.2. Informative References

- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC8193] Burbridge, T., Eardley, P., Bagnulo, M., and J. Schoenwaelder, "Information Model for Large-Scale Measurement Platforms (LMAPs)", RFC 8193, DOI 10.17487/RFC8193, August 2017, <<https://www.rfc-editor.org/info/rfc8193>>.
- [TR-181] Broadband Forum, "Device Data Model", <<http://cwmp-data-models.broadband-forum.org/>>.

Authors' Addresses

Barbara Stark
AT&T
Atlanta, GA,
United States of America

Email: barbara.stark@att.com

Maresh Jethanandani
VMware
California
United States of America

Email: mjethanandani@gmail.com

Network Working Group
Internet-Draft
Obsoletes: 6126,7557 (if approved)
Intended status: Standards Track
Expires: February 26, 2021

J. Chroboczek
IRIF, University of Paris-Diderot
D. Schinazi
Google LLC
August 25, 2020

The Babel Routing Protocol
draft-ietf-babel-rfc6126bis-20

Abstract

Babel is a loop-avoiding distance-vector routing protocol that is robust and efficient both in ordinary wired networks and in wireless mesh networks. This document describes the Babel routing protocol, and obsoletes RFCs 6126 and 7557.

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 February 26, 2021.

Copyright Notice

Copyright (c) 2020 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	3
1.1. Features	3
1.2. Limitations	4
1.3. Specification of Requirements	5
2. Conceptual Description of the Protocol	5
2.1. Costs, Metrics and Neighbourship	5
2.2. The Bellman-Ford Algorithm	6
2.3. Transient Loops in Bellman-Ford	6
2.4. Feasibility Conditions	7
2.5. Solving Starvation: Sequencing Routes	8
2.6. Requests	10
2.7. Multiple Routers	11
2.8. Overlapping Prefixes	12
3. Protocol Operation	12
3.1. Message Transmission and Reception	12
3.2. Data Structures	13
3.3. Acknowledgments and acknowledgment requests	17
3.4. Neighbour Acquisition	18
3.5. Routing Table Maintenance	21
3.6. Route Selection	25
3.7. Sending Updates	26
3.8. Explicit Requests	28
4. Protocol Encoding	32
4.1. Data Types	32
4.2. Packet Format	33
4.3. TLV Format	34
4.4. Sub-TLV Format	35
4.5. Parser state and encoding of updates	36
4.6. Details of Specific TLVs	37
4.7. Details of specific sub-TLVs	48
5. IANA Considerations	49
6. Security Considerations	52
7. Acknowledgments	53
8. References	53
8.1. Normative References	53
8.2. Informative References	54
Appendix A. Cost and Metric Computation	56
A.1. Maintaining Hello History	56
A.2. Cost Computation	57
A.3. Route selection and hysteresis	59
Appendix B. Protocol parameters	60
Appendix C. Route filtering	61
Appendix D. Considerations for protocol extensions	61
Appendix E. Stub Implementations	63
Appendix F. Compatibility with previous versions	64
Appendix G. Changes from previous versions	65

G.1.	Changes since RFC 6126	65
G.2.	Changes since draft-ietf-babel-rfc6126bis-00	65
G.3.	Changes since draft-ietf-babel-rfc6126bis-01	65
G.4.	Changes since draft-ietf-babel-rfc6126bis-02	66
G.5.	Changes since draft-ietf-babel-rfc6126bis-03	66
G.6.	Changes since draft-ietf-babel-rfc6126bis-03	67
G.7.	Changes since draft-ietf-babel-rfc6126bis-04	67
G.8.	Changes since draft-ietf-babel-rfc6126bis-05	67
G.9.	Changes since draft-ietf-babel-rfc6126bis-06	67
G.10.	Changes since draft-ietf-babel-rfc6126bis-07	67
G.11.	Changes since draft-ietf-babel-rfc6126bis-08	67
G.12.	Changes since draft-ietf-babel-rfc6126bis-09	68
G.13.	Changes since draft-ietf-babel-rfc6126bis-10	68
G.14.	Changes since draft-ietf-babel-rfc6126bis-11	68
G.15.	Changes since draft-ietf-babel-rfc6126bis-12	68
G.16.	Changes since draft-ietf-babel-rfc6126bis-13	69
G.17.	Changes since draft-ietf-babel-rfc6126bis-14	69
G.18.	Changes since draft-ietf-babel-rfc6126bis-15	69
G.19.	Changes since draft-ietf-babel-rfc6126bis-16	69
G.20.	Changes since draft-ietf-babel-rfc6126bis-17	69
G.21.	Changes since draft-ietf-babel-rfc6126bis-18	70
G.22.	Changes since draft-ietf-babel-rfc6126bis-19	70
Authors' Addresses		70

1. Introduction

Babel is a loop-avoiding distance-vector routing protocol that is designed to be robust and efficient both in networks using prefix-based routing and in networks using flat routing ("mesh networks"), and both in relatively stable wired networks and in highly dynamic wireless networks. This document describes the Babel routing protocol, and obsoletes [RFC6126] and [RFC7557].

1.1. Features

The main property that makes Babel suitable for unstable networks is that, unlike naive distance-vector routing protocols [RIP], it strongly limits the frequency and duration of routing pathologies such as routing loops and black-holes during reconvergence. Even after a mobility event is detected, a Babel network usually remains loop-free. Babel then quickly reconverges to a configuration that preserves the loop-freedom and connectedness of the network, but is not necessarily optimal; in many cases, this operation requires no packet exchanges at all. Babel then slowly converges, in a time on the scale of minutes, to an optimal configuration. This is achieved by using sequenced routes, a technique pioneered by Destination-Sequenced Distance-Vector routing [DSDV].

More precisely, Babel has the following properties:

- o when every prefix is originated by at most one router, Babel never suffers from routing loops;
- o when a single prefix is originated by multiple routers, Babel may occasionally create a transient routing loop for this particular prefix; this loop disappears in time proportional to the loop's diameter, and never again (up to an arbitrary garbage-collection (GC) time) will the routers involved participate in a routing loop for the same prefix;
- o assuming bounded packet loss rates, any routing black-holes that may appear after a mobility event are corrected in a time at most proportional to the network's diameter.

Babel has provisions for link quality estimation and for fairly arbitrary metrics. When configured suitably, Babel can implement shortest-path routing, or it may use a metric based, for example, on measured packet loss.

Babel nodes will successfully establish an association even when they are configured with different parameters. For example, a mobile node that is low on battery may choose to use larger time constants (hello and update intervals, etc.) than a node that has access to wall power. Conversely, a node that detects high levels of mobility may choose to use smaller time constants. The ability to build such heterogeneous networks makes Babel particularly adapted to the unmanaged or wireless environment.

Finally, Babel is a hybrid routing protocol, in the sense that it can carry routes for multiple network-layer protocols (IPv4 and IPv6), regardless of which protocol the Babel packets are themselves being carried over.

1.2. Limitations

Babel has two limitations that make it unsuitable for use in some environments. First, Babel relies on periodic routing table updates rather than using a reliable transport; hence, in large, stable networks it generates more traffic than protocols that only send updates when the network topology changes. In such networks, protocols such as OSPF [OSPF], IS-IS [IS-IS], or the Enhanced Interior Gateway Routing Protocol (EIGRP) [EIGRP] might be more suitable.

Second, unless the second algorithm described in Section 3.5.4 is implemented, Babel does impose a hold time when a prefix is

retracted. While this hold time does not apply to the exact prefix being retracted, and hence does not prevent fast reconvergence should it become available again, it does apply to any shorter prefix that covers it. This may make those implementations of Babel that do not implement the optional algorithm described in Section 3.5.4 unsuitable for use in networks that implement automatic prefix aggregation.

1.3. Specification of Requirements

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. Conceptual Description of the Protocol

Babel is a loop-avoiding distance vector protocol: it is based on the Bellman-Ford algorithm, just like the venerable RIP [RIP], but includes a number of refinements that either prevent loop formation altogether, or ensure that a loop disappears in a timely manner and doesn't form again.

Conceptually, Bellman-Ford is executed in parallel for every source of routing information (destination of data traffic). In the following discussion, we fix a source *S*; the reader will recall that the same algorithm is executed for all sources.

2.1. Costs, Metrics and Neighbourship

For every pair of neighbouring nodes *A* and *B*, Babel computes an abstract value known as the cost of the link from *A* to *B*, written *C*(*A*, *B*). Given a route between any two (not necessarily neighbouring) nodes, the metric of the route is the sum of the costs of all the links along the route. The goal of the routing algorithm is to compute, for every source *S*, the tree of routes of lowest metric to *S*.

Costs and metrics need not be integers. In general, they can be values in any algebra that satisfies two fairly general conditions (Section 3.5.2).

A Babel node periodically sends Hello messages to all of its neighbours; it also periodically sends an IHU ("I Heard You") message to every neighbour from which it has recently heard a Hello. From the information derived from Hello and IHU messages received from its

neighbour B, a node A computes the cost $C(A, B)$ of the link from A to B.

2.2. The Bellman-Ford Algorithm

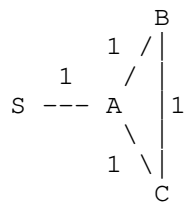
Every node A maintains two pieces of data: its estimated distance to S, written $D(A)$, and its next-hop router to S, written $NH(A)$. Initially, $D(S) = 0$, $D(A)$ is infinite, and $NH(A)$ is undefined.

Periodically, every node B sends to all of its neighbours a route update, a message containing $D(B)$. When a neighbour A of B receives the route update, it checks whether B is its selected next hop; if that is the case, then $NH(A)$ is set to B, and $D(A)$ is set to $C(A, B) + D(B)$. If that is not the case, then A compares $C(A, B) + D(B)$ to its current value of $D(A)$. If that value is smaller, meaning that the received update advertises a route that is better than the currently selected route, then $NH(A)$ is set to B, and $D(A)$ is set to $C(A, B) + D(B)$.

A number of refinements to this algorithm are possible, and are used by Babel. In particular, convergence speed may be increased by sending unscheduled "triggered updates" whenever a major change in the topology is detected, in addition to the regular, scheduled updates. Additionally, a node may maintain a number of alternate routes, which are being advertised by neighbours other than its selected neighbour, and which can be used immediately if the selected route were to fail.

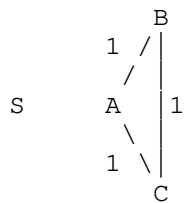
2.3. Transient Loops in Bellman-Ford

It is well known that a naive application of Bellman-Ford to distributed routing can cause transient loops after a topology change. Consider for example the following topology:



After convergence, $D(B) = D(C) = 2$, with $NH(B) = NH(C) = A$.

Suppose now that the link between S and A fails:



When it detects the failure of the link, A switches its next hop to B (which is still advertising a route to S with metric 2), and advertises a metric equal to 3, and then advertises a new route with metric 3. This process of nodes changing selected neighbours and increasing their metric continues until the advertised metric reaches "infinity", a value larger than all the metrics that the routing protocol is able to carry.

2.4. Feasibility Conditions

Bellman-Ford is a very robust algorithm: its convergence properties are preserved when routers delay route acquisition or when they discard some updates. Babel routers discard received route announcements unless they can prove that accepting them cannot possibly cause a routing loop.

More formally, we define a condition over route announcements, known as the "feasibility condition", that guarantees the absence of routing loops whenever all routers ignore route updates that do not satisfy the feasibility condition. In effect, this makes Bellman-Ford into a family of routing algorithms, parameterised by the feasibility condition.

Many different feasibility conditions are possible. For example, BGP can be modelled as being a distance-vector protocol with a (rather drastic) feasibility condition: a routing update is only accepted when the receiving node's AS number is not included in the update's AS-Path attribute (note that BGP's feasibility condition does not ensure the absence of transient "micro-loops" during reconvergence).

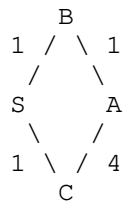
Another simple feasibility condition, used in the Destination-Sequenced Distance-Vector (DSDV) routing protocol [DSDV] and in the Ad hoc On-Demand Distance Vector (AODV) protocol [RFC3561], stems from the following observation: a routing loop can only arise after a router has switched to a route with a larger metric than the route that it had previously selected. Hence, one may define that a route is feasible when its metric at the local node would be no larger than the metric of the currently selected route, i.e., an announcement carrying a metric $D(B)$ is accepted by A when $C(A, B) + D(B) \leq D(A)$. If all routers obey this constraint, then the metric at every router

is nonincreasing, and the following invariant is always preserved: if A has selected B as its next hop, then $D(B) < D(A)$, which implies that the forwarding graph is loop-free.

Babel uses a slightly more refined feasibility condition, derived from EIGRP [DUAL]. Given a router A, define the feasibility distance of A, written $FD(A)$, as the smallest metric that A has ever advertised for S to any of its neighbours. An update sent by a neighbour B of A is feasible when the metric $D(B)$ advertised by B is strictly smaller than A's feasibility distance, i.e., when $D(B) < FD(A)$.

It is easy to see that this latter condition is no more restrictive than DSDV-feasibility. Suppose that node A obeys DSDV-feasibility; then $D(A)$ is nonincreasing, hence at all times $D(A) \leq FD(A)$. Suppose now that A receives a DSDV-feasible update that advertises a metric $D(B)$. Since the update is DSDV-feasible, $C(A, B) + D(B) \leq D(A)$, hence $D(B) < D(A)$, and since $D(A) \leq FD(A)$, $D(B) < FD(A)$.

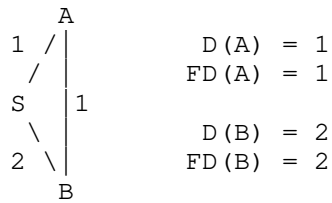
To see that it is strictly less restrictive, consider the following diagram, where A has selected the route through B, and $D(A) = FD(A) = 2$. Since $D(C) = 1 < FD(A)$, the alternate route through C is feasible for A, although its metric $C(A, C) + D(C) = 5$ is larger than that of the currently selected route:



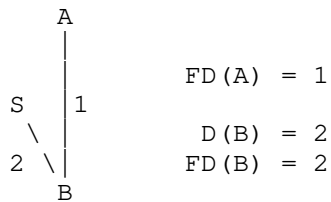
To show that this feasibility condition still guarantees loop-freedom, recall that at the time when A accepts an update from B, the metric $D(B)$ announced by B is no smaller than $FD(B)$; since it is smaller than $FD(A)$, at that point in time $FD(B) < FD(A)$. Since this property is preserved both when A sends updates and when it picks a different next hop, it remains true at all times, which ensures that the forwarding graph has no loops.

2.5. Solving Starvation: Sequencing Routes

Obviously, the feasibility conditions defined above cause starvation when a router runs out of feasible routes. Consider the following diagram, where both A and B have selected the direct route to S:



Suppose now that the link between A and S breaks:

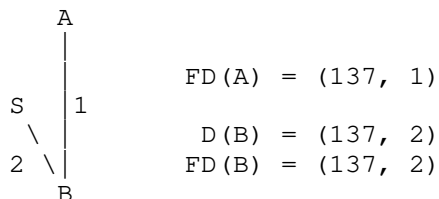


The only route available from A to S, the one that goes through B, is not feasible: A suffers from spurious starvation. At that point, the whole subtree suffering from starvation must be reset, which is essentially what EIGRP does when it performs a global synchronisation of all the routers in the starving subtree (the "active" phase of EIGRP).

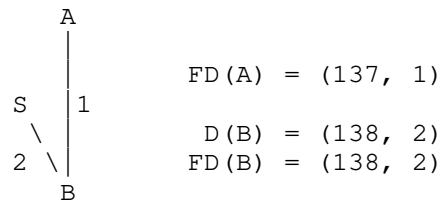
Babel reacts to starvation in a less drastic manner, by using sequenced routes, a technique introduced by DSDV and adopted by AODV. In addition to a metric, every route carries a sequence number, a nondecreasing integer that is propagated unchanged through the network and is only ever incremented by the source; a pair (s, m) , where s is a sequence number and m a metric, is called a distance.

A received update is feasible when either it is more recent than the feasibility distance maintained by the receiving node, or it is equally recent and the metric is strictly smaller. More formally, if $FD(A) = (s, m)$, then an update carrying the distance (s', m') is feasible when either $s' > s$, or $s = s'$ and $m' < m$.

Assuming the sequence number of S is 137, the diagram above becomes:



After S increases its sequence number, and the new sequence number is propagated to B, we have:



at which point the route through B becomes feasible again.

Note that while sequence numbers are used for determining feasibility, they are not used in route selection: a node ignores the sequence number when selecting the best route to a given destination (Section 3.6). Doing otherwise would cause route oscillation while a sequence number propagates through the network, and might even cause persistent blackholes with some exotic metrics.

2.6. Requests

In DSDV, the sequence number of a source is increased periodically. A route becomes feasible again after the source increases its sequence number, and the new sequence number is propagated through the network, which may, in general, require a significant amount of time.

Babel takes a different approach. When a node detects that it is suffering from a potentially spurious starvation, it sends an explicit request to the source for a new sequence number. This request is forwarded hop by hop to the source, with no regard to the feasibility condition. Upon receiving the request, the source increases its sequence number and broadcasts an update, which is forwarded to the requesting node.

Note that after a change in network topology not all such requests will, in general, reach the source, as some will be sent over links that are now broken. However, if the network is still connected, then at least one among the nodes suffering from spurious starvation has an (unfeasible) route to the source; hence, in the absence of packet loss, at least one such request will reach the source. (Resending requests a small number of times compensates for packet loss.)

Since requests are forwarded with no regard to the feasibility condition, they may, in general, be caught in a forwarding loop; this

is avoided by having nodes perform duplicate detection for the requests that they forward.

2.7. Multiple Routers

The above discussion assumes that each prefix is originated by a single router. In real networks, however, it is often necessary to have a single prefix originated by multiple routers: for example, the default route will be originated by all of the edge routers of a routing domain.

Since synchronising sequence numbers between distinct routers is problematic, Babel treats routes for the same prefix as distinct entities when they are originated by different routers: every route announcement carries the router-id of its originating router, and feasibility distances are not maintained per prefix, but per source, where a source is a pair of a router-id and a prefix. In effect, Babel guarantees loop-freedom for the forwarding graph to every source; since the union of multiple acyclic graphs is not in general acyclic, Babel does not in general guarantee loop-freedom when a prefix is originated by multiple routers, but any loops will be broken in a time at most proportional to the diameter of the loop -- as soon as an update has "gone around" the routing loop.

Consider for example the following topology, where A has selected the default route through S, and B has selected the one through S':

```

      1      1      1
:::/0 -- S --- A --- B --- S' -- :::/0

```

Suppose that both default routes fail at the same time; then nothing prevents A from switching to B, and B simultaneously switching to A. However, as soon as A has successfully advertised the new route to B, the route through A will become unfeasible for B. Conversely, as soon as B will have advertised the route through A, the route through B will become unfeasible for A.

In effect, the routing loop disappears at the latest when routing information has gone around the loop. Since this process can be delayed by lost packets, Babel makes certain efforts to ensure that updates are sent reliably after a router-id change (Section 3.7.2).

Additionally, after the routers have advertised the two routes, both sources will be in their source tables, which will prevent them from ever again participating in a routing loop involving routes from S and S' (up to the source GC time, which, available memory permitting, can be set to arbitrarily large values).

2.8. Overlapping Prefixes

In the above discussion, we have assumed that all prefixes are disjoint, as is the case in flat ("mesh") routing. In practice, however, prefixes may overlap: for example, the default route overlaps with all of the routes present in the network.

After a route fails, it is not correct in general to switch to a route that subsumes the failed route. Consider for example the following configuration:

```
          1      1
::/0 -- A --- B --- C
```

Suppose that node C fails. If B forwards packets destined to C by following the default route, a routing loop will form, and persist until A learns of B's retraction of the direct route to C. B avoids this pitfall by installing an "unreachable" route after a route is retracted; this route is maintained until it can be guaranteed that the former route has been retracted by all of B's neighbours (Section 3.5.4).

3. Protocol Operation

Every Babel speaker is assigned a router-id, which is an arbitrary string of 8 octets that is assumed unique across the routing domain. For example, router-ids could be assigned randomly, or they could be derived from a link-layer address. (The protocol encoding is slightly more compact when router-ids are assigned in the same manner as the IPv6 layer assigns host IDs; see the definition of the "R" flag in Section 4.6.9.)

3.1. Message Transmission and Reception

Babel protocol packets are sent in the body of a UDP datagram (as described in Section 4 below). Each Babel packet consists of zero or more TLVs. Most TLVs may contain sub-TLVs.

The protocol's control traffic can be carried indifferently over IPv6 or over IPv4, and prefixes of either address family can be announced over either protocol. Thus, there are at least two natural deployment models: using IPv6 exclusively for all control traffic, or running two distinct protocol instances, one for each address family. The exclusive use of IPv6 for all control traffic is RECOMMENDED, since using both protocols at the same time doubles the amount of traffic devoted to neighbour discovery and link quality estimation.

The source address of a Babel packet is always a unicast address, link-local in the case of IPv6. Babel packets may be sent to a well-known (link-local) multicast address or to a (link-local) unicast address. In normal operation, a Babel speaker sends both multicast and unicast packets to its neighbours.

With the exception of acknowledgments, all Babel TLVs can be sent to either unicast or multicast addresses, and their semantics does not depend on whether the destination is a unicast or a multicast address. Hence, a Babel speaker does not need to determine the destination address of a packet that it receives in order to interpret it.

A moderate amount of jitter may be applied to packets sent by a Babel speaker: outgoing TLVs are buffered and SHOULD be sent with a random delay. This is done for two purposes: it avoids synchronisation of multiple Babel speakers across a network [JITTER], and it allows for the aggregation of multiple TLVs into a single packet.

The maximum amount of delay a TLV can be subjected to depends on the TLV. When the protocol description specifies that a TLV is urgent (as in Section 3.7.2 and Section 3.8.2), then the TLV MUST be sent within a short time known as the urgent timeout (see Appendix B for recommended values). When the TLV is governed by a timeout explicitly included in a previous TLV, such as in the case of Acknowledgements (Section 4.6.4), Updates (Section 3.7) and IHUs (Section 3.4.2), then the TLV MUST be sent early enough to meet the explicit deadline (with a small margin to allow for propagation delays). In all other cases, the TLV SHOULD be sent out within one-half of the Multicast Hello interval.

In order to avoid packet drops (either at the sender or at the receiver), a delay SHOULD be introduced between successive packets sent out on the same interface, within the constraints of the previous paragraph. Note however that such packet pacing might impair the ability of some link layers (e.g., IEEE 802.11 [IEEE802.11]) to perform packet aggregation.

3.2. Data Structures

In this section, we give a description of the data structures that every Babel speaker maintains. This description is conceptual: a Babel speaker may use different data structures as long as the resulting protocol is the same as the one described in this document. For example, rather than maintaining a single table containing both selected and unselected (fallback) routes, as described in Section 3.2.6 below, an actual implementation would probably use two tables, one with selected routes and one with fallback routes.

3.2.1. Sequence number arithmetic

Sequence numbers (seqnos) appear in a number of Babel data structures, and they are interpreted as integers modulo 2^{16} . For the purposes of this document, arithmetic on sequence numbers is defined as follows.

Given a seqno s and a non-negative integer n , the sum of s and n is defined by

$$s + n \text{ (modulo } 2^{16}) = (s + n) \text{ MOD } 2^{16}$$

or, equivalently,

$$s + n \text{ (modulo } 2^{16}) = (s + n) \text{ AND } 65535$$

where MOD is the modulo operation yielding a non-negative integer and AND is the bitwise conjunction operation.

Given two sequence numbers s and s' , the relation s is less than s' ($s < s'$) is defined by

$$s < s' \text{ (modulo } 2^{16}) \text{ when } 0 < ((s' - s) \text{ MOD } 2^{16}) < 32768$$

or equivalently

$$s < s' \text{ (modulo } 2^{16}) \text{ when } s \neq s' \text{ and } ((s' - s) \text{ AND } 32768) = 0.$$

3.2.2. Node Sequence Number

A node's sequence number is a 16-bit integer that is included in route updates sent for routes originated by this node.

A node increments its sequence number (modulo 2^{16}) whenever it receives a request for a new sequence number (Section 3.8.1.2). A node SHOULD NOT increment its sequence number (seqno) spontaneously, since increasing seqnos makes it less likely that other nodes will have feasible alternate routes when their selected routes fail.

3.2.3. The Interface Table

The interface table contains the list of interfaces on which the node speaks the Babel protocol. Every interface table entry contains the interface's outgoing Multicast Hello seqno, a 16-bit integer that is sent with each Multicast Hello TLV on this interface and is incremented (modulo 2^{16}) whenever a Multicast Hello is sent. (Note that an interface's Multicast Hello seqno is unrelated to the node's seqno.)

There are two timers associated with each interface table entry. The periodic Multicast Hello timer governs the sending of scheduled Multicast Hello and IHU packets (Section 3.4). The periodic Update timer governs the sending of periodic route updates (Section 3.7.1). See Appendix B for suggested time constants.

3.2.4. The Neighbour Table

The neighbour table contains the list of all neighbouring interfaces from which a Babel packet has been recently received. The neighbour table is indexed by pairs of the form (interface, address), and every neighbour table entry contains the following data:

- o the local node's interface over which this neighbour is reachable;
- o the address of the neighbouring interface;
- o a history of recently received Multicast Hello packets from this neighbour; this can, for example, be a sequence of n bits, for some small value n , indicating which of the n hellos most recently sent by this neighbour have been received by the local node;
- o a history of recently received Unicast Hello packets from this neighbour;
- o the "transmission cost" value from the last IHU packet received from this neighbour, or FFFF hexadecimal (infinity) if the IHU hold timer for this neighbour has expired;
- o the expected incoming Multicast Hello sequence number for this neighbour, an integer modulo 2^{16} .
- o the expected incoming Unicast Hello sequence number for this neighbour, an integer modulo 2^{16} .
- o the outgoing Unicast Hello sequence number for this neighbour, an integer modulo 2^{16} that is sent with each Unicast Hello TLV to this neighbour and is incremented (modulo 2^{16}) whenever a Unicast Hello is sent. (Note that the outgoing Unicast Hello seqno for a neighbour is distinct from the interface's outgoing Multicast Hello seqno.)

There are three timers associated with each neighbour entry -- the multicast hello timer, which is set to the interval value carried by scheduled Multicast Hello TLVs sent by this neighbour, the unicast hello timer, which is set to the interval value carried by scheduled Unicast Hello TLVs, and the IHU timer, which is set to a small

multiple of the interval carried in IHU TLVs (see "IHU Hold time" in Appendix B for suggested values).

Note that the neighbour table is indexed by IP addresses, not by router-ids: neighbourship is a relationship between interfaces, not between nodes. Therefore, two nodes with multiple interfaces can participate in multiple neighbourship relationships, a situation that can notably arise when wireless nodes with multiple radios are involved.

3.2.5. The Source Table

The source table is used to record feasibility distances. It is indexed by triples of the form (prefix, plen, router-id), and every source table entry contains the following data:

- o the prefix (prefix, plen), where plen is the prefix length in bits, that this entry applies to;
- o the router-id of a router originating this prefix;
- o a pair (seqno, metric), this source's feasibility distance.

There is one timer associated with each entry in the source table -- the source garbage-collection timer. It is initialised to a time on the order of minutes and reset as specified in Section 3.7.3.

3.2.6. The Route Table

The route table contains the routes known to this node. It is indexed by triples of the form (prefix, plen, neighbour), and every route table entry contains the following data:

- o the source (prefix, plen, router-id) for which this route is advertised;
- o the neighbour (an entry in the neighbour table) that advertised this route;
- o the metric with which this route was advertised by the neighbour, or FFFF hexadecimal (infinity) for a recently retracted route;
- o the sequence number with which this route was advertised;
- o the next-hop address of this route;

- o a boolean flag indicating whether this route is selected, i.e., whether it is currently being used for forwarding and is being advertised.

There is one timer associated with each route table entry -- the route expiry timer. It is initialised and reset as specified in Section 3.5.3.

Note that there are two distinct (seqno, metric) pairs associated to each route: the route's distance, which is stored in the route table, and the feasibility distance, stored in the source table and shared between all routes with the same source.

3.2.7. The Table of Pending Seqno Requests

The table of pending seqno requests contains a list of seqno requests that the local node has sent (either because they have been originated locally, or because they were forwarded) and to which no reply has been received yet. This table is indexed by triples of the form (prefix, plen, router-id), and every entry in this table contains the following data:

- o the prefix, plen, router-id, and seqno being requested;
- o the neighbour, if any, on behalf of which we are forwarding this request;
- o a small integer indicating the number of times that this request will be resent if it remains unsatisfied.

There is one timer associated with each pending seqno request; it governs both the resending of requests and their expiry.

3.3. Acknowledgments and acknowledgment requests

A Babel speaker may request that a neighbour receiving a given packet reply with an explicit acknowledgment within a given time. While the use of acknowledgment requests is optional, every Babel speaker MUST be able to reply to such a request.

An acknowledgment MUST be sent to a unicast destination. On the other hand, acknowledgment requests may be sent to either unicast or multicast destinations, in which case they request an acknowledgment from all of the receiving nodes.

When to request acknowledgments is a matter of local policy; the simplest strategy is to never request acknowledgments and to rely on periodic updates to ensure that any reachable routes are eventually

propagated throughout the routing domain. In order to improve convergence speed and reduce the amount of control traffic, acknowledgment requests MAY be used in order to reliably send urgent updates (Section 3.7.2) and retractions (Section 3.5.4), especially when the number of neighbours on a given interface is small. Since Babel is designed to deal gracefully with packet loss on unreliable media, sending all packets with acknowledgment requests is not necessary, and NOT RECOMMENDED, as the acknowledgments cause additional traffic and may force additional Address Resolution Protocol (ARP) or Neighbour Discovery (ND) exchanges.

3.4. Neighbour Acquisition

Neighbour acquisition is the process by which a Babel node discovers the set of neighbours heard over each of its interfaces and ascertains bidirectional reachability. On unreliable media, neighbour acquisition additionally provides some statistics that may be useful for link quality computation.

Before it can exchange routing information with a neighbour, a Babel node MUST create an entry for that neighbour in the neighbour table. When to do that is implementation-specific; suitable strategies include creating an entry when any Babel packet is received, or creating an entry when a Hello TLV is parsed. Similarly, in order to conserve system resources, an implementation SHOULD discard an entry when it has been unused for long enough; suitable strategies include dropping the neighbour after a timeout, and dropping a neighbour when the associated Hello histories become empty (see Appendix A.2).

3.4.1. Reverse Reachability Detection

Every Babel node sends Hello TLVs to its neighbours to indicate that it is alive, at regular or irregular intervals. Each Hello TLV carries an increasing (modulo 2^{16}) sequence number and an upper bound on the time interval until the next Hello of the same type (see below). If the time interval is set to 0, then the Hello TLV does not establish a new promise: the deadline carried by the previous Hello of the same type still applies to the next Hello (if the most recent scheduled Hello of the right kind was received at time t_0 and carried interval i , then the previous promise of sending another Hello before time $t_0 + i$ still holds). We say that a Hello is "scheduled" if it carries a non-zero interval, and "unscheduled" otherwise.

There are two kinds of Hellos: Multicast Hellos, which use a per-interface Hello counter (the Multicast Hello seqno), and Unicast Hellos, which use a per-neighbour counter (the Unicast Hello seqno). A Multicast Hello with a given seqno MUST be sent to all neighbours

on a given interface, either by sending it to a multicast address or by sending it to one unicast address per neighbour (hence, the term "Multicast Hello" is a slight misnomer). A Unicast Hello carrying a given seqno should normally be sent to just one neighbour (over unicast), since the sequence numbers of different neighbours are not in general synchronised.

Multicast Hellos sent over multicast can be used for neighbour discovery; hence, a node SHOULD send periodic (scheduled) Multicast Hellos unless neighbour discovery is performed by means outside of the Babel protocol. A node MAY send Unicast Hellos or unscheduled Hellos of either kind for any reason, such as reducing the amount of multicast traffic or improving reliability on link technologies with poor support for link-layer multicast.

A node MAY send a scheduled Hello ahead of time. A node MAY change its scheduled Hello interval. The Hello interval MAY be decreased at any time; it MAY be increased immediately before sending a Hello TLV, but SHOULD NOT be increased at other times. (Equivalently, a node SHOULD send a scheduled Hello immediately after increasing its Hello interval.)

How to deal with received Hello TLVs and what statistics to maintain are considered local implementation matters; typically, a node will maintain some sort of history of recently received Hellos. An example of a suitable algorithm is described in Appendix A.1.

After receiving a Hello, or determining that it has missed one, the node recomputes the association's cost (Section 3.4.3) and runs the route selection procedure (Section 3.6).

3.4.2. Bidirectional Reachability Detection

In order to establish bidirectional reachability, every node sends periodic IHU ("I Heard You") TLVs to each of its neighbours. Since IHUs carry an explicit interval value, they MAY be sent less often than Hellos in order to reduce the amount of routing traffic in dense networks; in particular, they SHOULD be sent less often than Hellos over links with little packet loss. While IHUs are conceptually unicast, they MAY be sent to a multicast address in order to avoid an ARP or Neighbour Discovery exchange and to aggregate multiple IHUs into a single packet.

In addition to the periodic IHUs, a node MAY, at any time, send an unscheduled IHU packet. It MAY also, at any time, decrease its IHU interval, and it MAY increase its IHU interval immediately before sending an IHU, but SHOULD NOT increase it at any other time.

(Equivalently, a node SHOULD send an extra IHU immediately after increasing its Hello interval.)

Every IHU TLV contains two pieces of data: the link's rxcost (reception cost) from the sender's perspective, used by the neighbour for computing link costs (Section 3.4.3), and the interval between periodic IHU packets. A node receiving an IHU sets the value of the txcost (transmission cost) maintained in the neighbour table to the value contained in the IHU, and resets the IHU timer associated to this neighbour to a small multiple of the interval value received in the IHU (see "IHU Hold time" in Appendix B for suggested values). When a neighbour's IHU timer expires, the neighbour's txcost is set to infinity.

After updating a neighbour's txcost, the receiving node recomputes the neighbour's cost (Section 3.4.3) and runs the route selection procedure (Section 3.6).

3.4.3. Cost Computation

A neighbourhood association's link cost is computed from the values maintained in the neighbour table: the statistics kept in the neighbour table about the reception of Hellos, and the txcost computed from received IHU packets.

For every neighbour, a Babel node computes a value known as this neighbour's rxcost. This value is usually derived from the Hello history, which may be combined with other data, such as statistics maintained by the link layer. The rxcost is sent to a neighbour in each IHU.

Since nodes do not necessarily send periodic Unicast Hellos but do usually send periodic Multicast Hellos (Section 3.4.1), a node SHOULD use an algorithm that yields a finite rxcost when only Multicast Hellos are received, unless interoperability with nodes that only send Multicast Hellos is not required.

How the txcost and rxcost are combined in order to compute a link's cost is a matter of local policy; as far as Babel's correctness is concerned, only the following conditions MUST be satisfied:

- o the cost is strictly positive;
- o if no Hello TLVs of either kind were received recently, then the cost is infinite;
- o if the txcost is infinite, then the cost is infinite.

See Appendix A.2 for RECOMMENDED strategies for computing a link's cost.

3.5. Routing Table Maintenance

Conceptually, a Babel update is a quintuple (prefix, plen, router-id, seqno, metric), where (prefix, plen) is the prefix for which a route is being advertised, router-id is the router-id of the router originating this update, seqno is a nondecreasing (modulo 2^{16}) integer that carries the originating router seqno, and metric is the announced metric.

Before being accepted, an update is checked against the feasibility condition (Section 3.5.1), which ensures that the route does not create a routing loop. If the feasibility condition is not satisfied, the update is either ignored or prevents the route from being selected, as described in Section 3.5.3. If the feasibility condition is satisfied, then the update cannot possibly cause a routing loop.

3.5.1. The Feasibility Condition

The feasibility condition is applied to all received updates. The feasibility condition compares the metric in the received update with the metrics of the updates previously sent by the receiving node; updates that fail the feasibility condition, and therefore have metrics large enough to cause a routing loop, are either ignored or prevent the resulting route from being selected.

A feasibility distance is a pair (seqno, metric), where seqno is an integer modulo 2^{16} and metric is a positive integer. Feasibility distances are compared lexicographically, with the first component inverted: we say that a distance (seqno, metric) is strictly better than a distance (seqno', metric'), written

$$(\text{seqno}, \text{metric}) < (\text{seqno}', \text{metric}')$$

when

$$\text{seqno} > \text{seqno}' \text{ or } (\text{seqno} = \text{seqno}' \text{ and } \text{metric} < \text{metric}')$$

where sequence numbers are compared modulo 2^{16} .

Given a source (prefix, plen, router-id), a node's feasibility distance for this source is the minimum, according to the ordering defined above, of the distances of all the finite updates ever sent by this particular node for the prefix (prefix, plen) and the given

router-id. Feasibility distances are maintained in the source table, the exact procedure is given in Section 3.7.3.

A received update is feasible when either it is a retraction (its metric is FFFF hexadecimal), or the advertised distance is strictly better, in the sense defined above, than the feasibility distance for the corresponding source. More precisely, a route advertisement carrying the quintuple (prefix, plen, router-id, seqno, metric) is feasible if one of the following conditions holds:

- o metric is infinite; or
- o no entry exists in the source table indexed by (prefix, plen, router-id); or
- o an entry (prefix, plen, router-id, seqno', metric') exists in the source table, and either
 - * seqno' < seqno or
 - * seqno = seqno' and metric < metric'.

Note that the feasibility condition considers the metric advertised by the neighbour, not the route's metric; hence, a fluctuation in a neighbour's cost cannot render a selected route unfeasible. Note further that retractions (updates with infinite metric) are always feasible, since they cannot possibly cause a routing loop.

3.5.2. Metric Computation

A route's metric is computed from the metric advertised by the neighbour and the neighbour's link cost. Just like cost computation, metric computation is considered a local policy matter; as far as Babel is concerned, the function $M(c, m)$ used for computing a metric from a locally computed link cost c and the metric m advertised by a neighbour MUST only satisfy the following conditions:

- o if c is infinite, then $M(c, m)$ is infinite;
- o M is strictly monotonic: $M(c, m) > m$.

Additionally, the metric SHOULD satisfy the following condition:

- o M is left-distributive: if $m \leq m'$, then $M(c, m) \leq M(c, m')$.

While strict monotonicity is essential to the integrity of the network (persistent routing loops may arise if it is not satisfied), left distributivity is not: if it is not satisfied, Babel will still

converge to a loop-free configuration, but might not reach a global optimum (in fact, a global optimum may not even exist).

The conditions above are easily satisfied by using the additive metric, i.e., by defining $M(c, m) = c + m$. Since the additive metric is useful with a large range of cost computation strategies, it is the RECOMMENDED default. See also Appendix C, which describes a technique that makes it possible to tweak the values of individual metrics without running the risk of creating routing loops.

3.5.3. Route Acquisition

When a Babel node receives an update (prefix, plen, router-id, seqno, metric) from a neighbour neigh, it checks whether it already has a route table entry indexed by (prefix, plen, neigh).

If no such entry exists:

- o if the update is unfeasible, it MAY be ignored;
- o if the metric is infinite (the update is a retraction of a route we do not know about), the update is ignored;
- o otherwise, a new entry is created in the route table, indexed by (prefix, plen, neigh), with source equal to (prefix, plen, router-id), seqno equal to seqno and an advertised metric equal to the metric carried by the update.

If such an entry exists:

- o if the entry is currently selected, the update is unfeasible, and the router-id of the update is equal to the router-id of the entry, then the update MAY be ignored;
- o otherwise, the entry's sequence number, advertised metric, metric, and router-id are updated and, if the advertised metric is not infinite, the route's expiry timer is reset to a small multiple of the Interval value included in the update (see "Route Hold time" in Appendix B for suggested values). If the update is unfeasible, then the (now unfeasible) entry MUST be immediately unselected. If the update caused the router-id of the entry to change, an update (possibly a retraction) MUST be sent in a timely manner as described in Section 3.7.2.

Note that the route table may contain unfeasible routes, either because they were created by an unfeasible update or due to a metric fluctuation. Such routes are never selected, since they are not known to be loop-free; should all the feasible routes become

unusable, however, the unfeasible routes can be made feasible and therefore possible to select by sending requests along them (see Section 3.8.2).

When a route's expiry timer triggers, the behaviour depends on whether the route's metric is finite. If the metric is finite, it is set to infinity and the expiry timer is reset. If the metric is already infinite, the route is flushed from the route table.

After the route table is updated, the route selection procedure (Section 3.6) is run.

3.5.4. Hold Time

When a prefix P is retracted, because all routes are unfeasible or have an infinite metric (whether due to the expiry timer or to other reasons), and a shorter prefix P' that covers P is reachable, P' cannot in general be used for routing packets destined to P without running the risk of creating a routing loop (Section 2.8).

To avoid this issue, whenever a prefix P is retracted, a route table entry with infinite metric is maintained as described in Section 3.5.3 above. As long as this entry is maintained, packets destined to an address within P MUST NOT be forwarded by following a route for a shorter prefix. This entry is removed as soon as a finite-metric update for prefix P is received and the resulting route selected. If no such update is forthcoming, the infinite metric entry SHOULD be maintained at least until it is guaranteed that no neighbour has selected the current node as next-hop for prefix P. This can be achieved by either:

- o waiting until the route's expiry timer has expired (Section 3.5.3);
- o sending a retraction with an acknowledgment request (Section 3.3) to every reachable neighbour that has not explicitly retracted prefix P, and waiting for all acknowledgments.

The former option is simpler and ensures that at that point, any routes for prefix P pointing at the current node have expired. However, since the expiry time can be as high as a few minutes, doing that prevents automatic aggregation by creating spurious black-holes for aggregated routes. The latter option is RECOMMENDED as it dramatically reduces the time for which a prefix is unreachable in the presence of aggregated routes.

3.6. Route Selection

Route selection is the process by which a single route for a given prefix is selected to be used for forwarding packets and to be re-advertised to a node's neighbours.

Babel is designed to allow flexible route selection policies. As far as the algorithm's correctness is concerned, the route selection policy **MUST** only satisfy the following properties:

- o a route with infinite metric (a retracted route) is never selected;
- o an unfeasible route is never selected.

Babel nodes using different route selection strategies will interoperate and not create routing loops as long as these two properties hold.

Route selection **MUST NOT** take seqnos into account: a route **MUST NOT** be preferred just because it carries a higher (more recent) seqno. Doing otherwise would cause route oscillation while a new seqno propagates across the network, and might create persistent blackholes if the metric being used is not left-distributive (Section 3.5.2).

The obvious route selection strategy is to pick, for every destination, the feasible route with minimal metric. When all metrics are stable, this approach ensures convergence to a tree of shortest paths (assuming that the metric is left-distributive, see Section 3.5.2). There are two reasons, however, why this strategy may lead to instability in the presence of continuously varying metrics. First, if two parallel routes oscillate around a common value, then the smallest metric strategy will keep switching between the two. Second, when a route is selected, congestion along it increases, which might increase packet loss, which in turn could cause its metric to increase; this is a feedback loop, of the kind that is prone to causing persistent oscillations.

In order to limit these kinds of instabilities, a route selection strategy **SHOULD** include some form of hysteresis, i.e., be sensitive to a route's history: if a route is currently selected, then the strategy should only switch to a different route if the latter has been consistently good for some period of time. A **RECOMMENDED** hysteresis algorithm is given in Appendix A.3.

After the route selection procedure is run, triggered updates (Section 3.7.2) and requests (Section 3.8.2) are sent.

3.7. Sending Updates

A Babel speaker advertises to its neighbours its set of selected routes. Normally, this is done by sending one or more multicast packets containing Update TLVs on all of its connected interfaces; however, on link technologies where multicast is significantly more expensive than unicast, a node MAY choose to send multiple copies of updates in unicast packets, especially when the number of neighbours is small.

Additionally, in order to ensure that any black-holes are reliably cleared in a timely manner, a Babel node may send retractions (updates with an infinite metric) for any recently retracted prefixes.

If an update is for a route injected into the Babel domain by the local node (e.g., it carries the address of a local interface, the prefix of a directly attached network, or a prefix redistributed from a different routing protocol), the router-id is set to the local node's router-id, the metric is set to some arbitrary finite value (typically 0), and the seqno is set to the local router's sequence number.

If an update is for a route learned from another Babel speaker, the router-id and sequence number are copied from the route table entry, and the metric is computed as specified in Section 3.5.2.

3.7.1. Periodic Updates

Every Babel speaker MUST advertise each of its selected routes on every interface, at least once every Update interval. Since Babel doesn't suffer from routing loops (there is no "counting to infinity") and relies heavily on triggered updates (Section 3.7.2), this full dump only needs to happen infrequently (see Appendix B for suggested intervals).

3.7.2. Triggered Updates

In addition to periodic routing updates, a Babel speaker sends unscheduled, or triggered, updates in order to inform its neighbours of a significant change in the network topology.

A change of router-id for the selected route to a given prefix may be indicative of a routing loop in formation; hence, whenever it changes the selected router-id for a given destination, a node MUST send an update as an urgent TLV (as defined in Section 3.1). Additionally, it SHOULD make a reasonable attempt at ensuring that all reachable neighbours receive this update.

There are two strategies for ensuring that. If the number of neighbours is small, then it is reasonable to send the update together with an acknowledgment request; the update is resent until all neighbours have acknowledged the packet, up to some number of times. If the number of neighbours is large, however, requesting acknowledgments from all of them might cause a non-negligible amount of network traffic; in that case, it may be preferable to simply repeat the update some reasonable number of times (say, 3 for wireless and 2 for wired links). The number of copies **MUST NOT** exceed 5, and the copies **SHOULD** be separated by a small delay, as described in Section 3.1.

A route retraction is somewhat less worrying: if the route retraction doesn't reach all neighbours, a black-hole might be created, which, unlike a routing loop, does not endanger the integrity of the network. When a route is retracted, a node **SHOULD** send a triggered update and **SHOULD** make a reasonable attempt at ensuring that all neighbours receive this retraction.

Finally, a node **MAY** send a triggered update when the metric for a given prefix changes in a significant manner, due to a received update, because a link's cost has changed, or because a different next hop has been selected. A node **SHOULD NOT** send triggered updates for other reasons, such as when there is a minor fluctuation in a route's metric, when the selected next hop changes without inducing a significant change to the route's metric, or to propagate a new sequence number (except to satisfy a request, as specified in Section 3.8).

3.7.3. Maintaining Feasibility Distances

Before sending an update (prefix, plen, router-id, seqno, metric) with finite metric (i.e., not a route retraction), a Babel node updates the feasibility distance maintained in the source table. This is done as follows.

If no entry indexed by (prefix, plen, router-id) exists in the source table, then one is created with value (prefix, plen, router-id, seqno, metric).

If an entry (prefix, plen, router-id, seqno', metric') exists, then it is updated as follows:

- o if seqno > seqno', then seqno' := seqno, metric' := metric;
- o if seqno = seqno' and metric' > metric, then metric' := metric;
- o otherwise, nothing needs to be done.

The garbage-collection timer for the entry is then reset. Note that the feasibility distance is not updated and the garbage-collection timer is not reset when a retraction (an update with infinite metric) is sent.

When the garbage-collection timer expires, the entry is removed from the source table.

3.7.4. Split Horizon

When running over a transitive, symmetric link technology, e.g., a point-to-point link or a wired LAN technology such as Ethernet, a Babel node SHOULD use an optimisation known as split horizon. When split horizon is used on a given interface, a routing update for prefix P is not sent on the particular interface over which the selected route towards prefix P was learnt.

Split horizon SHOULD NOT be applied to an interface unless the interface is known to be symmetric and transitive; in particular, split horizon is not applicable to decentralised wireless link technologies (e.g., IEEE 802.11 in ad hoc mode) when routing updates are sent over multicast.

3.8. Explicit Requests

In normal operation, a node's route table is populated by the regular and triggered updates sent by its neighbours. Under some circumstances, however, a node sends explicit requests in order to cause a resynchronisation with the source after a mobility event or to prevent a route from spuriously expiring.

The Babel protocol provides two kinds of explicit requests: route requests, which simply request an update for a given prefix, and seqno requests, which request an update for a given prefix with a specific sequence number. The former are never forwarded; the latter are forwarded if they cannot be satisfied by the receiver.

3.8.1. Handling Requests

Upon receiving a request, a node either forwards the request or sends an update in reply to the request, as described in the following sections. If this causes an update to be sent, the update is either sent to a multicast address on the interface on which the request was received, or to the unicast address of the neighbour that sent the request.

The exact behaviour is different for route requests and seqno requests.

3.8.1.1. Route Requests

When a node receives a route request for a given prefix, it checks its route table for a selected route to this exact prefix. If such a route exists, it **MUST** send an update (over unicast or over multicast); if such a route does not exist, it **MUST** send a retraction for that prefix.

When a node receives a wildcard route request, it **SHOULD** send a full route table dump. Full route dumps **SHOULD** be rate-limited, especially if they are sent over multicast.

3.8.1.2. Seqno Requests

When a node receives a seqno request for a given router-id and sequence number, it checks whether its route table contains a selected entry for that prefix. If a selected route for the given prefix exists, it has finite metric, and either the router-ids are different or the router-ids are equal and the entry's sequence number is no smaller (modulo 2^{16}) than the requested sequence number, the node **MUST** send an update for the given prefix. If the router-ids match but the requested seqno is larger (modulo 2^{16}) than the route entry's, the node compares the router-id against its own router-id. If the router-id is its own, then it increases its sequence number by 1 (modulo 2^{16}) and sends an update. A node **MUST NOT** increase its sequence number by more than 1 in reaction to a single seqno request.

Otherwise, if the requested router-id is not its own, the received node consults the hop count field of the request. If the hop count is 2 or more, and the node is advertising the prefix to its neighbours, the node selects a neighbour to forward the request to as follows:

- o if the node has one or more feasible routes toward the requested prefix with a next hop that is not the requesting node, then the node **MUST** forward the request to the next hop of one such route;
- o otherwise, if the node has one or more (not feasible) routes to the requested prefix with a next hop that is not the requesting node, then the node **SHOULD** forward the request to the next hop of one such route.

In order to actually forward the request, the node decrements the hop count and sends the request in a unicast packet destined to the selected neighbour. The forwarded request **SHOULD** be sent as an urgent TLV (as defined in Section 3.1).

A node SHOULD maintain a list of recently forwarded seqno requests and forward the reply (an update with a seqno sufficiently large to satisfy the request) as an urgent TLV (as defined in Section 3.1). A node SHOULD compare every incoming seqno request against its list of recently forwarded seqno requests and avoid forwarding it if it is redundant (i.e., if it has recently sent a request with the same prefix, router-id and a seqno that is not smaller modulo 2^{16}).

Since the request-forwarding mechanism does not necessarily obey the feasibility condition, it may get caught in routing loops; hence, requests carry a hop count to limit the time during which they remain in the network. However, since requests are only ever forwarded as unicast packets, the initial hop count need not be kept particularly low, and performing an expanding horizon search is not necessary. A single request MUST NOT be duplicated: it MUST NOT be forwarded to a multicast address, and it MUST NOT be forwarded to multiple neighbours. However, if a seqno request is resent by its originator, the subsequent copies may be forwarded to a different neighbour than the initial one.

3.8.2. Sending Requests

A Babel node MAY send a route or seqno request at any time, to a multicast or a unicast address; there is only one case when originating requests is required (Section 3.8.2.1).

3.8.2.1. Avoiding Starvation

When a route is retracted or expires, a Babel node usually switches to another feasible route for the same prefix. It may be the case, however, that no such routes are available.

A node that has lost all feasible routes to a given destination but still has unexpired unfeasible routes to that destination MUST send a seqno request; if it doesn't have any such routes, it MAY still send a seqno request. The router-id of the request is set to the router-id of the route that it has just lost, and the requested seqno is the value contained in the source table plus 1. The request carries a hop count, which is used as a last-resort mechanism to ensure that it eventually vanishes from the network; it MAY be set to any value that is larger than the diameter of the network (64 is a suitable default value).

If the node has any (unfeasible) routes to the requested destination, then it MUST send the request to at least one of the next-hop neighbours that advertised these routes, and SHOULD send it to all of them; in any case, it MAY send the request to any other neighbours, whether they advertise a route to the requested destination or not.

A simple implementation strategy is therefore to unconditionally multicast the request over all interfaces.

Similar requests will be sent by other nodes that are affected by the route's loss. If the network is still connected, and assuming no packet loss, then at least one of these requests will be forwarded to the source, resulting in a route being advertised with a new sequence number. (Due to duplicate suppression, only a small number of such requests are expected to actually reach the source.)

In order to compensate for packet loss, a node SHOULD repeat such a request a small number of times if no route becomes feasible within a short time (see "Request Timeout" in Appendix B for suggested values). In the presence of heavy packet loss, however, all such requests might be lost; in that case, the mechanism in the next section will eventually ensure that a new seqno is received.

3.8.2.2. Dealing with Unfeasible Updates

When a route's metric increases, a node might receive an unfeasible update for a route that it has currently selected. As specified in Section 3.5.1, the receiving node will either ignore the update or unselect the route.

In order to keep routes from spuriously expiring because they have become unfeasible, a node SHOULD send a unicast seqno request when it receives an unfeasible update for a route that is currently selected. The requested sequence number is computed from the source table as in Section 3.8.2.1 above.

Additionally, since metric computation does not necessarily coincide with the delay in propagating updates, a node might receive an unfeasible update from a currently unselected neighbour that is preferable to the currently selected route (e.g., because it has a much smaller metric); in that case, the node SHOULD send a unicast seqno request to the neighbour that advertised the preferable update.

3.8.2.3. Preventing Routes from Expiring

In normal operation, a route's expiry timer never triggers: since a route's hold time is computed from an explicit interval included in Update TLVs, a new update (possibly a retraction) should arrive in time to prevent a route from expiring.

In the presence of packet loss, however, it may be the case that no update is successfully received for an extended period of time, causing a route to expire. In order to avoid such spurious expiry, shortly before a selected route expires, a Babel node SHOULD send a

unicast route request to the neighbour that advertised this route; since nodes always send either updates or retractions in response to non-wildcard route requests (Section 3.8.1.1), this will usually result in the route being either refreshed or retracted.

4. Protocol Encoding

A Babel packet MUST be sent as the body of a UDP datagram, with network-layer hop count set to 1, destined to a well-known multicast address or to a unicast address, over IPv4 or IPv6; in the case of IPv6, these addresses are link-local. Both the source and destination UDP port are set to a well-known port number. A Babel packet MUST be silently ignored unless its source address is either a link-local IPv6 address or an IPv4 address belonging to the local network, and its source port is the well-known Babel port. It MAY be silently ignored if its destination address is a global IPv6 address.

In order to minimise the number of packets being sent while avoiding lower-layer fragmentation, a Babel node SHOULD maximise the size of the packets it sends, up to the outgoing interface's MTU adjusted for lower-layer headers (28 octets for UDP over IPv4, 48 octets for UDP over IPv6). It MUST NOT send packets larger than the attached interface's MTU adjusted for lower-layer headers or 512 octets, whichever is larger, but not exceeding $2^{16} - 1$ adjusted for lower-layer headers. Every Babel speaker MUST be able to receive packets that are as large as any attached interface's MTU adjusted for lower-layer headers or 512 octets, whichever is larger. Babel packets MUST NOT be sent in IPv6 Jumbograms [RFC2675].

4.1. Data Types

4.1.1. Representation of integers

All multi-octet fields that represent integers are encoded with the most significant octet first (in Big-Endian format [IEN137], also called Network Order). The base protocol only carries unsigned values; if an extension needs to carry signed values, it will need to specify their encoding (e.g., two's complement).

4.1.2. Interval

Relative times are carried as 16-bit values specifying a number of centiseconds (hundredths of a second). This allows times up to roughly 11 minutes with a granularity of 10ms, which should cover all reasonable applications of Babel (see also Appendix B).

4.1.3. Router-Id

A router-id is an arbitrary 8-octet value. A router-id MUST NOT consist of either all binary zeroes (0000000000000000 hexadecimal) or all binary ones (FFFFFFFFFFFFFFFF hexadecimal).

4.1.4. Address

Since the bulk of the protocol is taken by addresses, multiple ways of encoding addresses are defined. Additionally, within Update TLVs a common subnet prefix may be omitted when multiple addresses are sent in a single packet -- this is known as address compression (Section 4.6.9).

Address encodings:

- o AE 0: wildcard address. The value is 0 octets long.
- o AE 1: IPv4 address. Compression is allowed. 4 octets or less.
- o AE 2: IPv6 address. Compression is allowed. 16 octets or less.
- o AE 3: link-local IPv6 address. Compression is not allowed. The value is 8 octets long, a prefix of fe80::/64 is implied.

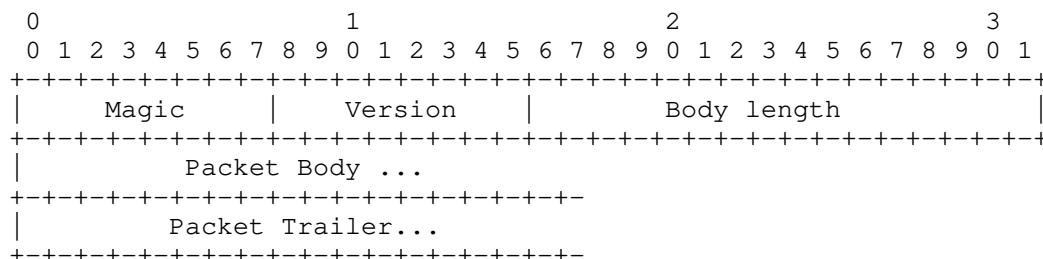
The address family associated to an address encoding is either IPv4 or IPv6; it is undefined for AE 0, IPv4 for AE 1, and IPv6 for AEs 2 and 3.

4.1.5. Prefixes

A network prefix is encoded just like a network address, but it is stored in the smallest number of octets that are enough to hold the significant bits (up to the prefix length).

4.2. Packet Format

A Babel packet consists of a 4-octet header, followed by a sequence of TLVs (the packet body), optionally followed by a second sequence of TLVs (the packet trailer). The format is designed to be extensible; see Appendix D for extensibility considerations.



Fields :

Magic The arbitrary but carefully chosen value 42 (decimal); packets with a first octet different from 42 MUST be silently ignored.

Version This document specifies version 2 of the Babel protocol. Packets with a second octet different from 2 MUST be silently ignored.

Body length The length in octets of the body following the packet header (excluding the Magic, Version and Body length fields, and excluding the packet trailer).

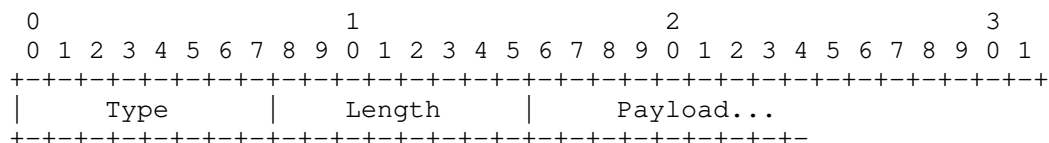
Packet Body The packet body; a sequence of TLVs.

Packet Trailer The packet trailer; another sequence of TLVs.

The packet body and trailer are both sequences of TLVs. The packet body is the normal place to store TLVs; the packet trailer only contains specialised TLVs that do not need to be protected by cryptographic security mechanisms. When parsing the trailer, the receiver MUST ignore any TLV unless its definition explicitly states that it is allowed to appear there. Among the TLVs defined in this document, only Pad1 and PadN are allowed in the trailer; since these TLVs are ignored in any case, an implementation MAY silently ignore the packet trailer without even parsing it, unless it implements at least one protocol extension that defines TLVs that are allowed to appear in the trailer.

4.3. TLV Format

With the exception of Pad1, all TLVs have the following structure:



Fields :

Type The type of the TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

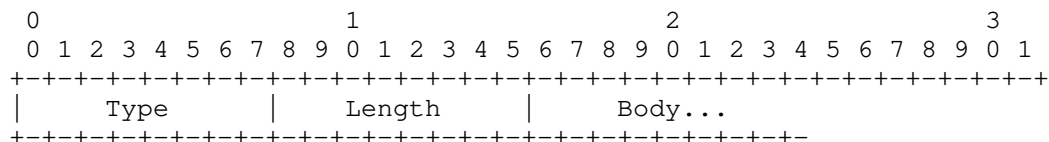
Payload The TLV payload, which consists of a body and, for selected TLV types, an optional list of sub-TLVs.

TLVs with an unknown type value MUST be silently ignored.

4.4. Sub-TLV Format

Every TLV carries an explicit length in its header; however, most TLVs are self-terminating, in the sense that it is possible to determine the length of the body without reference to the explicit Length field. If a TLV has a self-terminating format, then the space between the natural size of the TLV and the size announced in the Length field may be used to store a sequence of sub-TLVs.

Sub-TLVs have the same structure as TLVs. With the exception of Pad1, all TLVs have the following structure:



Fields :

Type The type of the sub-TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

Body The sub-TLV body, the interpretation of which depends on both the type of the sub-TLV and the type of the TLV within which it is embedded.

The most-significant bit of the sub-TLV type (the bit with value 80 hexadecimal), is called the mandatory bit; in other words, sub-TLV types 128 through 255 have the mandatory bit set. This bit indicates how to handle unknown sub-TLVs. If the mandatory bit is not set, then an unknown sub-TLV MUST be silently ignored, and the rest of the TLV is processed normally. If the mandatory bit is set, then the whole enclosing TLV MUST be silently ignored (except for updating the parser state by a Router-Id, Next-Hop or Update TLV, as described in the next section).

4.5. Parser state and encoding of updates

In a large network, the bulk of Babel traffic consists of route updates; hence, some care has been given to encoding them efficiently. The data conceptually contained in an update (Section 3.5) is split into three pieces:

- o the prefix, seqno and metric are contained in the Update TLV itself (Section 4.6.9);
- o the router-id is taken from Router-Id TLV that precedes the Update TLV, and may be shared among multiple Update TLVs (Section 4.6.7);
- o the next hop is taken either from the source-address of the network-layer packet that contains the Babel packet, or from an explicit Next-Hop TLV (Section 4.6.8).

In addition to the above, an Update TLV can omit a prefix of the prefix being announced, which is then extracted from the preceding Update TLV in the same address family (IPv4 or IPv6). Finally, as a special optimisation for the case when a router-id coincides with the interface-id part of an IPv6 address, Router-ID TLV itself may be omitted and the router-id derived from the low-order bits of the advertised prefix (Section 4.6.9).

In order to implement these compression techniques, Babel uses a stateful parser: a TLV may refer to data from a previous TLV. The parser state consists of the following pieces of data:

- o for each address encoding that allows compression, the current default prefix; this is undefined at the start of the packet, and is updated by each Update TLV with the Prefix flag set (Section 4.6.9);
- o for each address family (IPv4 or IPv6), the current next-hop; this is the source address of the enclosing packet for the matching address family at the start of a packet, and is updated by each Next-Hop TLV (Section 4.6.8);

- o the current router-id; this is undefined at the start of the packet, and is updated by each Router-ID TLV (Section 4.6.7) and by each Update TLV with Router-Id flag set.

Since the parser state must be identical across implementations, it is updated before checking for mandatory sub-TLVs: parsing a TLV MUST update the parser state even if the TLV is otherwise ignored due to an unknown mandatory sub-TLV or for any other reason.

None of the TLVs that modify the parser state are allowed in the packet trailer; hence, an implementation may choose to use a dedicated stateless parser to parse the packet trailer.

4.6. Details of Specific TLVs

4.6.1. Pad1

```

0
0 1 2 3 4 5 6 7
+---+---+---+---+
|   Type = 0   |
+---+---+---+---+
```

Fields :

Type Set to 0 to indicate a Pad1 TLV.

This TLV is silently ignored on reception. It is allowed in the packet trailer.

4.6.2. PadN

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type = 1   | Length | MBZ... |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

Fields :

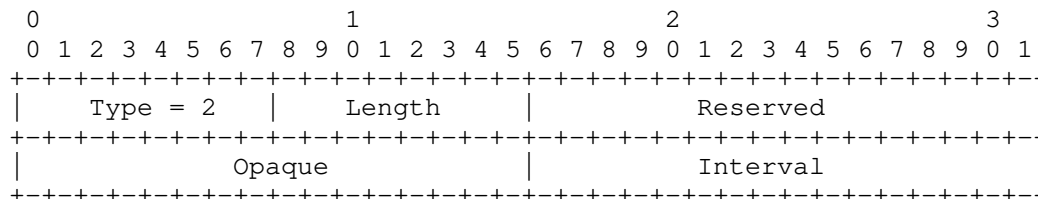
Type Set to 1 to indicate a PadN TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

MBZ Must be zero, set to 0 on transmission.

This TLV is silently ignored on reception. It is allowed in the packet trailer.

4.6.3. Acknowledgment Request



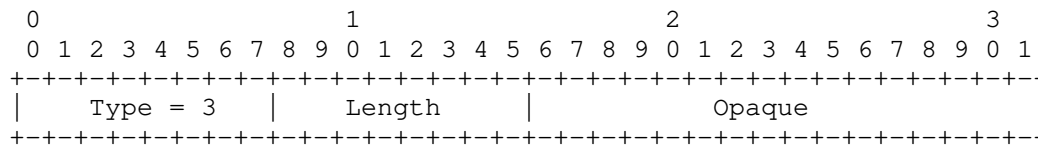
This TLV requests that the receiver send an Acknowledgment TLV within the number of centiseconds specified by the Interval field.

Fields :

- Type Set to 2 to indicate an Acknowledgment Request TLV.
- Length The length of the body in octets, exclusive of the Type and Length fields.
- Reserved Sent as 0 and MUST be ignored on reception.
- Opaque An arbitrary value that will be echoed in the receiver's Acknowledgment TLV.
- Interval A time interval in centiseconds after which the sender will assume that this packet has been lost. This MUST NOT be 0. The receiver MUST send an Acknowledgment TLV before this time has elapsed (with a margin allowing for propagation time).

This TLV is self-terminating, and allows sub-TLVs.

4.6.4. Acknowledgment



This TLV is sent by a node upon receiving an Acknowledgment Request.

Fields :

Type Set to 3 to indicate an Acknowledgment TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

Opaque Set to the Opaque value of the Acknowledgment Request that prompted this Acknowledgment.

Since Opaque values are not globally unique, this TLV MUST be sent to a unicast address.

This TLV is self-terminating, and allows sub-TLVs.

4.6.5. Hello

0										1										2										3									
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1								
Type = 4										Length										Flags																			
Seqno										Interval																													

This TLV is used for neighbour discovery and for determining a neighbour's reception cost.

Fields :

Type Set to 4 to indicate a Hello TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

Flags The individual bits of this field specify special handling of this TLV (see below).

Seqno If the Unicast flag is set, this is the value of the sending node's outgoing Unicast Hello seqno for this neighbour. Otherwise, it is the sending node's outgoing Multicast Hello seqno for this interface.

Interval If non-zero, this is an upper bound, expressed in centiseconds, on the time after which the sending node will send a new scheduled Hello TLV with the same setting of the Unicast flag. If this is 0, then this Hello represents an unscheduled Hello, and doesn't carry any new information about times at which Hellos are sent.

The Flags field is interpreted as follows:

```

      0                               1
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
+---+---+---+---+---+---+---+---+
|U|X|X|X|X|X|X|X|X|X|X|X|X|X|X|X|
+---+---+---+---+---+---+---+---+

```

- o U (Unicast) flag (8000 hexadecimal): if set, then this Hello represents a Unicast Hello, otherwise it represents a Multicast Hello;
- o X: all other bits MUST be sent as 0 and silently ignored on reception.

Every time a Hello is sent, the corresponding seqno counter MUST be incremented. Since there is a single seqno counter for all the Multicast Hellos sent by a given node over a given interface, if the Unicast flag is not set, this TLV MUST be sent to all neighbors on this link, which can be achieved by sending to a multicast destination, or by sending multiple packets to the unicast addresses of all reachable neighbours. Conversely, if the Unicast flag is set, this TLV MUST be sent to a single neighbour, which can be achieved by sending to a unicast destination. In order to avoid large discontinuities in link quality, multiple Hello TLVs SHOULD NOT be sent in the same packet.

This TLV is self-terminating, and allows sub-TLVs.

4.6.6. IHU

```

      0                               1                               2                               3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type = 5   |   Length   |   AE   |   Reserved   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           Rxcost           |           Interval           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           Address...           |
+---+---+---+---+---+---+---+---+

```

An IHU ("I Heard You") TLV is used for confirming bidirectional reachability and carrying a link's transmission cost.

Fields :

Type Set to 5 to indicate an IHU TLV.

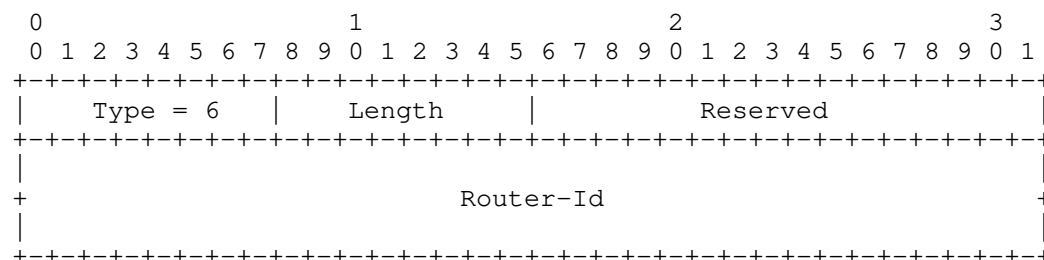
Length	The length of the body in octets, exclusive of the Type and Length fields.
AE	The encoding of the Address field. This should be 1 or 3 in most cases. As an optimisation, it MAY be 0 if the TLV is sent to a unicast address, if the association is over a point-to-point link, or when bidirectional reachability is ascertained by means outside of the Babel protocol.
Reserved	Sent as 0 and MUST be ignored on reception.
Rxcost	The rxcost according to the sending node of the interface whose address is specified in the Address field. The value FFFF hexadecimal (infinity) indicates that this interface is unreachable.
Interval	An upper bound, expressed in centiseconds, on the time after which the sending node will send a new IHU; this MUST NOT be 0. The receiving node will use this value in order to compute a hold time for this symmetric association.
Address	The address of the destination node, in the format specified by the AE field. Address compression is not allowed.

Conceptually, an IHU is destined to a single neighbour. However, IHU TLVs contain an explicit destination address, and MAY be sent to a multicast address, as this allows aggregation of IHUs destined to distinct neighbours into a single packet and avoids the need for an ARP or Neighbour Discovery exchange when a neighbour is not being used for data traffic.

IHU TLVs with an unknown value in the AE field MUST be silently ignored.

This TLV is self-terminating, and allows sub-TLVs.

4.6.7. Router-Id



A Router-Id TLV establishes a router-id that is implied by subsequent Update TLVs, as described in Section 4.5. This TLV sets the router-id even if it is otherwise ignored due to an unknown mandatory sub-TLV.

Fields :

Type Set to 6 to indicate a Router-Id TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

Reserved Sent as 0 and MUST be ignored on reception.

Router-Id The router-id for routes advertised in subsequent Update TLVs. This MUST NOT consist of all zeroes or all ones.

This TLV is self-terminating, and allows sub-TLVs.

4.6.8. Next Hop

```

      0               1               2               3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-----+-----+-----+-----+-----+-----+-----+-----+
|   Type = 7   |   Length   |   AE   |   Reserved   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|   Next hop...   |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

A Next Hop TLV establishes a next-hop address for a given address family (IPv4 or IPv6) that is implied in subsequent Update TLVs, as described in Section 4.5. This TLV sets up the next-hop for subsequent Update TLVs even if it is otherwise ignored due to an unknown mandatory sub-TLV.

Fields :

Type Set to 7 to indicate a Next Hop TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

AE The encoding of the Address field. This SHOULD be 1 (IPv4) or 3 (link-local IPv6), and MUST NOT be 0.

Reserved Sent as 0 and MUST be ignored on reception.

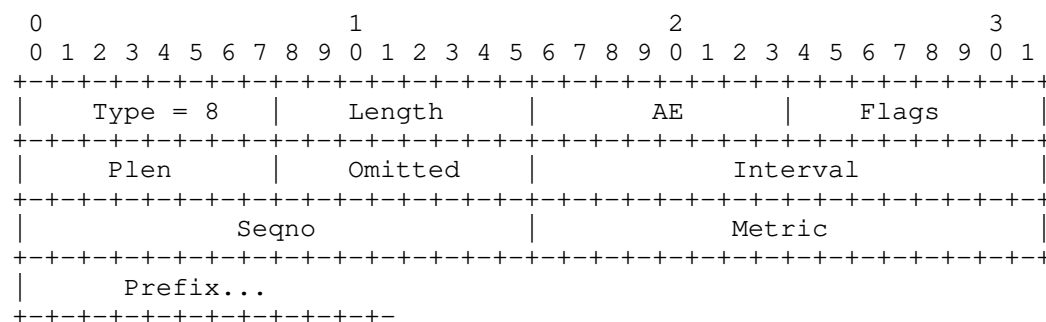
Next hop The next-hop address advertised by subsequent Update TLVs, for this address family.

When the address family matches the network-layer protocol that this packet is transported over, a Next Hop TLV is not needed: in the absence of a Next Hop TLV in a given address family, the next hop address is taken to be the source address of the packet.

Next Hop TLVs with an unknown value for the AE field MUST be silently ignored.

This TLV is self-terminating, and allows sub-TLVs.

4.6.9. Update



An Update TLV advertises or retracts a route. As an optimisation, it can optionally have the side effect of establishing a new implied router-id and a new default prefix, as described in Section 4.5.

Fields :

Type	Set to 8 to indicate an Update TLV.
Length	The length of the body in octets, exclusive of the Type and Length fields.
AE	The encoding of the Prefix field.
Flags	The individual bits of this field specify special handling of this TLV (see below).
Plen	The length in bits of the advertised prefix. If AE is 3 (link-local IPv6), Omitted MUST be 0.
Omitted	The number of octets that have been omitted at the beginning of the advertised prefix and that should be taken

from a preceding Update TLV in the same address family with the Prefix flag set.

Interval An upper bound, expressed in centiseconds, on the time after which the sending node will send a new update for this prefix. This MUST NOT be 0. The receiving node will use this value to compute a hold time for the route table entry. The value FFFF hexadecimal (infinity) expresses that this announcement will not be repeated unless a request is received (Section 3.8.2.3).

Seqno The originator's sequence number for this update.

Metric The sender's metric for this route. The value FFFF hexadecimal (infinity) means that this is a route retraction.

Prefix The prefix being advertised. This field's size is (Plen/8 - Omitted) rounded upwards.

The Flags field is interpreted as follows:

```

 0 1 2 3 4 5 6 7
+-+--+--+--+--+--+
|P|R|X|X|X|X|X|X|
+-+--+--+--+--+--+

```

- o P (Prefix) flag (80 hexadecimal): if set, then this Update establishes a new default prefix for subsequent Update TLVs with a matching address encoding within the same packet, even if this TLV is otherwise ignored due to an unknown mandatory sub-TLV;
- o R (Router-Id) flag (40 hexadecimal): if set, then this TLV establishes a new default router-id for this TLV and subsequent Update TLVs in the same packet, even if this TLV is otherwise ignored due to an unknown mandatory sub-TLV. This router-id is computed from the first address of the advertised prefix as follows:
 - * if the length of the address is 8 octets or more, then the new router-id is taken from the 8 last octets of the address;
 - * if the length of the address is smaller than 8 octets, then the new router-id consists of the required number of zero octets followed by the address, i.e., the address is stored on the right of the router-id. For example, for an IPv4 address, the router-id consists of 4 octets of zeroes followed by the IPv4 address.

- o X: all other bits MUST be sent as 0 and silently ignored on reception.

The prefix being advertised by an Update TLV is computed as follows:

- o the first Omitted octets of the prefix are taken from the previous Update TLV with the Prefix flag set and the same address encoding, even if it was ignored due to an unknown mandatory sub-TLV; if Omitted is not zero and there is no such TLV, then this Update MUST be ignored;
- o the next $(\text{Plen}/8 - \text{Omitted})$ rounded upwards octets are taken from the Prefix field;
- o if Plen is not a multiple of 8, then any bits beyond Plen (i.e., the low-order $(8 - \text{Plen} \bmod 8)$ bits of the last octet) are cleared;
- o the remaining octets are set to 0.

If the Metric field is finite, the router-id of the originating node for this announcement is taken from the prefix advertised by this Update if the Router-Id flag is set, computed as described above. Otherwise, it is taken either from the preceding Router-Id TLV, or the preceding Update TLV with the Router-Id flag set, whichever comes last, even if that TLV is otherwise ignored due to an unknown mandatory sub-TLV; if there is no suitable TLV, then this update is ignored.

The next-hop address for this update is taken from the last preceding Next Hop TLV with a matching address family (IPv4 or IPv6) in the same packet even if it was otherwise ignored due to an unknown mandatory sub-TLV; if no such TLV exists, it is taken from the network-layer source address of this packet if it belongs to the same address family as the prefix being announced; otherwise, this Update MUST be ignored.

If the metric field is FFFF hexadecimal, this TLV specifies a retraction. In that case, the router-id, next-hop and seqno are not used. AE MAY then be 0, in which case this Update retracts all of the routes previously advertised by the sending interface. If the metric is finite, AE MUST NOT be 0; Update TLVs with finite metric and AE equal to 0 MUST be ignored. If the metric is infinite and AE is 0, Plen and Omitted MUST both be 0; Update TLVs that do not satisfy this requirement MUST be ignored.

Update TLVs with an unknown value in the AE field MUST be silently ignored.

This TLV is self-terminating, and allows sub-TLVs.

4.6.10. Route Request

```

      0               1               2               3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-----+-----+-----+-----+-----+-----+-----+-----+
|  Type = 9  |  Length  |    AE    |    Plen   |
+-----+-----+-----+-----+-----+-----+-----+-----+
| Prefix...  |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

A Route Request TLV prompts the receiver to send an update for a given prefix, or a full route table dump.

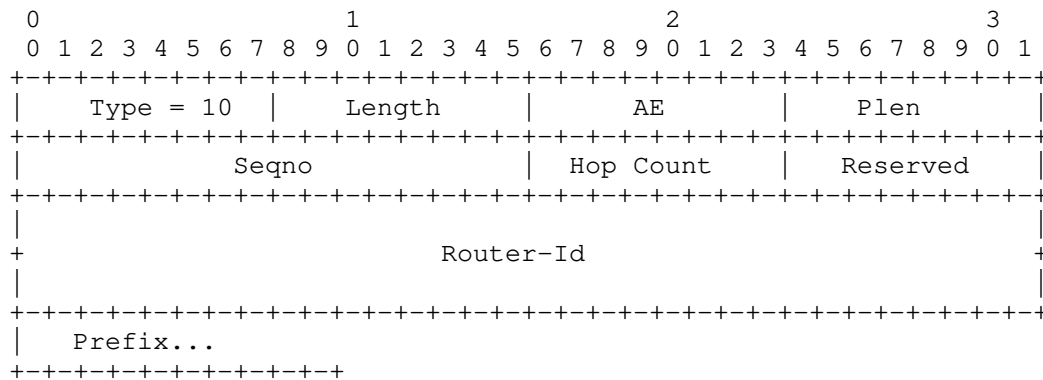
Fields :

Type	Set to 9 to indicate a Route Request TLV.
Length	The length of the body in octets, exclusive of the Type and Length fields.
AE	The encoding of the Prefix field. The value 0 specifies that this is a request for a full route table dump (a wildcard request).
Plen	The length in bits of the requested prefix.
Prefix	The prefix being requested. This field's size is Plen/8 rounded upwards.

A Request TLV prompts the receiver to send an update message (possibly a retraction) for the prefix specified by the AE, Plen, and Prefix fields, or a full dump of its route table if AE is 0 (in which case Plen must be 0 and Prefix is of length 0). A Request TLV with AE set to 0 and Plen not set to 0 MUST be ignored.

This TLV is self-terminating, and allows sub-TLVs.

4.6.11. Seqno Request



A Seqno Request TLV prompts the receiver to send an Update for a given prefix with a given sequence number, or to forward the request further if it cannot be satisfied locally.

Fields :

Type	Set to 10 to indicate a Seqno Request TLV.
Length	The length of the body in octets, exclusive of the Type and Length fields.
AE	The encoding of the Prefix field. This MUST NOT be 0.
Plen	The length in bits of the requested prefix.
Seqno	The sequence number that is being requested.
Hop Count	The maximum number of times that this TLV may be forwarded, plus 1. This MUST NOT be 0.
Reserved	Sent as 0 and MUST be ignored on reception.
Router-Id	The Router-Id that is being requested. This MUST NOT consist of all zeroes or all ones.
Prefix	The prefix being requested. This field's size is Plen/8 rounded upwards.

A Seqno Request TLV prompts the receiving node to send a finite-metric Update for the prefix specified by the AE, Plen, and Prefix fields, with either a router-id different from what is specified by the Router-Id field, or a Seqno no less (modulo 2^{16}) than what is specified by the Seqno field. If this request cannot be satisfied

locally, then it is forwarded according to the rules set out in Section 3.8.1.2.

While a Seqno Request MAY be sent to a multicast address, it MUST NOT be forwarded to a multicast address and MUST NOT be forwarded to more than one neighbour. A request MUST NOT be forwarded if its Hop Count field is 1.

This TLV is self-terminating, and allows sub-TLVs.

4.7. Details of specific sub-TLVs

4.7.1. Pad1

```

  0 1 2 3 4 5 6 7
+---+---+---+---+
|   Type = 0   |
+---+---+---+---+
```

Fields :

Type Set to 0 to indicate a Pad1 sub-TLV.

This sub-TLV is silently ignored on reception. It is allowed within any TLV that allows sub-TLVs.

4.7.2. PadN

```

      0                      1                      2                      3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type = 1   | Length | MBZ...
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

Fields :

Type Set to 1 to indicate a PadN sub-TLV.

Length The length of the body in octets, exclusive of the Type and Length fields.

MBZ Must be zero, set to 0 on transmission.

This sub-TLV is silently ignored on reception. It is allowed within any TLV that allows sub-TLVs.

5. IANA Considerations

IANA has registered the UDP port number 6696, called "babel", for use by the Babel protocol.

IANA has registered the IPv6 multicast group ff02::1:6 and the IPv4 multicast group 224.0.0.111 for use by the Babel protocol.

IANA has created a registry called "Babel TLV Types". The allocation policy for this registry is Specification Required [RFC8126] for Types 0-223, and Experimental Use for Types 224-254. The values in this registry are as follows:

Type	Name	Reference
0	Pad1	this document
1	PadN	this document
2	Acknowledgment Request	this document
3	Acknowledgment	this document
4	Hello	this document
5	IHU	this document
6	Router-Id	this document
7	Next Hop	this document
8	Update	this document
9	Route Request	this document
10	Seqno Request	this document
11	TS/PC	[RFC7298]
12	HMAC	[RFC7298]
13	Source-specific Update	[BABEL-SS]
14	Source-specific Request	[BABEL-SS]
15	Source-specific Seqno Request	[BABEL-SS]

16	MAC	[BABEL-MAC]
17	PC	[BABEL-MAC]
18	Challenge Request	[BABEL-MAC]
19	Challenge Reply	[BABEL-MAC]
20-223	Unassigned	
224-254	Reserved for Experimental Use	this document
255	Reserved for expansion of the type space	this document

IANA has created a registry called "Babel sub-TLV Types". The allocation policy for this registry is Specification Required for Types 0-111 and 128-239, and Experimental Use for Types 112-126 and 240-254. The values in this registry are as follows:

Type	Name	Reference
0	Pad1	this document
1	PadN	this document
2	Diversity	[BABEL-DIVERSITY]
3	Timestamp	[BABEL-RTT]
4-111	Unassigned	
112-126	Reserved for Experimental Use	this document
127	Reserved for expansion of the type space	this document
128	Source Prefix	[BABEL-SS]
129-239	Unassigned	
240-254	Reserved for Experimental Use	this document
255	Reserved for expansion of the type space	this document

IANA is instructed to create a registry called "Babel Address Encodings". The allocation policy for this registry is Specification Required for Address Encodings (AEs) 0-223, and Experimental Use for AEs 224-254. The values in this registry are as follows:

AE	Name	Reference
0	Wildcard address	this document
1	IPv4 address	this document
2	IPv6 address	this document
3	Link-local IPv6 address	this document
4-223	Unassigned	
224-254	Reserved for Experimental Use	this document
255	Reserved for expansion of the AE space	this document

IANA has created a registry called "Babel Flags Values". The allocation policy for this registry is Specification Required. IANA is instructed to rename this registry to "Babel Update Flags Values". The values in this registry are as follows:

Bit	Name	Reference
0	Default prefix	this document
1	Default Router-ID	this document
2-7	Unassigned	

IANA is instructed to create a new registry called "Babel Hello Flags Values". The allocation policy for this registry is Specification Required. The initial values in this registry are as follows:

Bit	Name	Reference
0	Unicast	this document
1-15	Unassigned	

IANA is instructed to replace all references to RFCs 6126 and 7557 in all of the registries mentioned above by references to this document.

6. Security Considerations

As defined in this document, Babel is a completely insecure protocol. Without additional security mechanisms, Babel trusts any information it receives in plaintext UDP datagrams and acts on it. An attacker that is present on the local network can impact Babel operation in a variety of ways; for example they can:

- o spoof a Babel packet, and redirect traffic by announcing a route with a smaller metric, a larger sequence number, or a longer prefix;
- o spoof a malformed packet, which could cause an insufficiently robust implementation to crash or interfere with the rest of the network;
- o replay a previously captured Babel packet, which could cause traffic to be redirected, blackholed or otherwise interfere with the network.

When carried over IPv6, Babel packets are ignored unless they are sent from a link-local IPv6 address; since routers don't forward link-local IPv6 packets, this mitigates the attacks outlined above by restricting them to on-link attackers. No such natural protection exists when Babel packets are carried over IPv4, which is one of the reasons why it is recommended to deploy Babel over IPv6 (Section 3.1).

It is usually difficult to ensure that packets arriving at a Babel node are trusted, even in the case where the local link is believed to be secure. For that reason, it is RECOMMENDED that all Babel traffic be protected by an application-layer cryptographic protocol. There are currently two suitable mechanisms, which implement different tradeoffs between implementation simplicity and security:

- o Babel over DTLS [BABEL-DTLS] runs the majority of Babel traffic over DTLS, and leverages DTLS to authenticate nodes and provide confidentiality and integrity protection;
- o MAC authentication [BABEL-MAC] appends a message authentication code (MAC) to every Babel packet to prove that it originated at a node that knows a shared secret, and includes sufficient additional information to prove that the packet is fresh (not replayed).

Both mechanisms enable nodes to ignore packets generated by attackers without the proper credentials. They also ensure integrity of messages and prevent message replay. While Babel-DTLS supports asymmetric keying and ensures confidentiality, Babel-MAC has a much more limited scope (see Sections 1.1, 1.2 and 7 of [BABEL-MAC]). Since Babel-MAC is simpler and more lightweight, it is recommended in preference to Babel-DTLS in deployments where its limitations are acceptable, i.e., when symmetric keying is sufficient and where the routing information is not considered confidential.

Every implementation of Babel SHOULD implement BABEL-MAC.

One should be aware that the information that a mobile Babel node announces to the whole routing domain is sufficient to determine the mobile node's physical location with reasonable precision, which might cause privacy concerns even if the control traffic is protected from unauthenticated attackers by a cryptographic mechanism such as Babel-DTLS. This issue may be mitigated somewhat by using randomly chosen router-ids and randomly chosen IP addresses, and changing them often enough.

7. Acknowledgments

A number of people have contributed text and ideas to this specification. The authors are particularly indebted to Matthieu Boutier, Gwendoline Chouasne, Margaret Cullen, Donald Eastlake, Toke Hoiland-Jorgensen, Benjamin Kaduk, Joao Sobrinho and Martin Vigoureux. Earlier versions of this document greatly benefited from the input of Joel Halpern. The address compression technique was inspired by [PACKETBB].

8. References

8.1. Normative References

- [BABEL-MAC] Do, C., Kolodziejek, W., and J. Chroboczek, "MAC authentication for the Babel routing protocol", Internet Draft draft-ietf-babel-hmac-10, August 2019.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997.
- [RFC793] Postel, J., "Transmission Control Protocol", RFC 793, DOI 10.17487/RFC0793, September 1981, <<https://www.rfc-editor.org/info/rfc793>>.
- [RFC8126] Cotton, M., Leiba, B., and T. Narten, "Guidelines for Writing an IANA Considerations Section in RFCs", BCP 26, RFC 8126, June 2017.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017.

8.2. Informative References

- [BABEL-DIVERSITY] Chroboczek, J., "Diversity Routing for the Babel Routing Protocol", draft-chroboczek-babel-diversity-routing-01 (work in progress), February 2016.
- [BABEL-DTLS] Decimo, A., Schinazi, D., and J. Chroboczek, "Babel Routing Protocol over Datagram Transport Layer Security", Internet Draft draft-ietf-babel-dtls-10, June 2020.
- [BABEL-RTT] Jonglez, B. and J. Chroboczek, "Delay-based Metric Extension for the Babel Routing Protocol", draft-ietf-babel-rtt-extension-00 (work in progress), April 2019.
- [BABEL-SS] Boutier, M. and J. Chroboczek, "Source-Specific Routing in Babel", draft-ietf-babel-source-specific-05 (work in progress), April 2019.
- [DSDV] Perkins, C. and P. Bhagwat, "Highly Dynamic Destination-Sequenced Distance-Vector Routing (DSDV) for Mobile Computers", ACM SIGCOMM'94 Conference on Communications Architectures, Protocols and Applications 234-244, 1994.

- [DUAL] Garcia Luna Aceves, J., "Loop-Free Routing Using Diffusing Computations", IEEE/ACM Transactions on Networking 1:1, February 1993.
- [EIGRP] Albrightson, B., Garcia Luna Aceves, J., and J. Boyle, "EIGRP -- a Fast Routing Protocol Based on Distance Vectors", Proc. Interop 94, 1994.
- [ETX] De Couto, D., Aguayo, D., Bicket, J., and R. Morris, "A high-throughput path metric for multi-hop wireless networks", Proc. MobiCom 2003, 2003.
- [IEEE802.11] IEEE, "IEEE Standard for Information technology-- Telecommunications and information exchange between systems Local and metropolitan area networks--Specific requirements Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications", IEEE 802.11-2012, DOI 10.1109/ieeestd.2012.6178212, April 2012.
- [IEN137] Cohen, D., "On holy wars and a plea for peace", IEN 137, April 1980.
- [IS-IS] Standardization, I. O. F., "Information technology -- Telecommunications and information exchange between systems -- Intermediate System to Intermediate System intra-domain routeing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode network service (ISO 8473)", ISO/IEC 10589:2002, 2002.
- [JITTER] Floyd, S. and V. Jacobson, "The synchronization of periodic routing messages", IEEE/ACM Transactions on Networking 2, 2, 122-136, April 1994.
- [OSPF] Moy, J., "OSPF Version 2", RFC 2328, April 1998.
- [PACKETBB] Clausen, T., Dearlove, C., Dean, J., and C. Adjih, "Generalized Mobile Ad Hoc Network (MANET) Packet/Message Format", RFC 5444, February 2009.
- [RFC2675] Borman, D., Deering, S., and R. Hinden, "IPv6 Jumbograms", RFC 2675, DOI 10.17487/RFC2675, August 1999.

- [RFC3561] Perkins, C., Belding-Royer, E., and S. Das, "Ad hoc On-Demand Distance Vector (AODV) Routing", RFC 3561, DOI 10.17487/RFC3561, July 2003, <<https://www.rfc-editor.org/info/rfc3561>>.
- [RFC6126] Chroboczek, J., "The Babel Routing Protocol", RFC 6126, DOI 10.17487/RFC6126, April 2011.
- [RFC7298] Ovsienko, D., "Babel Hashed Message Authentication Code (HMAC) Cryptographic Authentication", RFC 7298, DOI 10.17487/RFC7298, July 2014.
- [RFC7557] Chroboczek, J., "Extension Mechanism for the Babel Routing Protocol", RFC 7557, DOI 10.17487/RFC7557, May 2015.
- [RIP] Malkin, G., "RIP Version 2", RFC 2453, November 1998.

Appendix A. Cost and Metric Computation

The strategy for computing link costs and route metrics is a local matter; Babel itself only requires that it comply with the conditions given in Section 3.4.3 and Section 3.5.2. Different nodes may use different strategies in a single network and may use different strategies on different interface types. This section describes a set of strategies that have been found to work well in actual networks.

In summary, a node maintains per-neighbour statistics about the last 16 received Hello TLVs of each kind (Appendix A.1), it computes costs by using the 2-out-of-3 strategy (Appendix A.2.1) on wired links, and ETX (Appendix A.2.2) on wireless links. It uses an additive algebra for metric computation (Section 3.5.2).

A.1. Maintaining Hello History

For each neighbour, a node maintains two sets of Hello history, one for each kind of Hello, and an expected sequence number, one for Multicast and one for Unicast Hellos. Each Hello history is a vector of 16 bits, where a 1 value represents a received Hello, and a 0 value a missed Hello. For each kind of Hello, the expected sequence number, written *ne*, is the sequence number that is expected to be carried by the next received Hello from this neighbour.

Whenever it receives a Hello packet of a given kind from a neighbour, a node compares the received sequence number *nr* for that kind of Hello with its expected sequence number *ne*. Depending on the outcome of this comparison, one of the following actions is taken:

- o if the two differ by more than 16 (modulo 2^{16}), then the sending node has probably rebooted and lost its sequence number; the whole associated neighbour table entry is flushed and a new one is created;
- o otherwise, if the received nr is smaller (modulo 2^{16}) than the expected sequence number ne, then the sending node has increased its Hello interval without us noticing; the receiving node removes the last (ne - nr) entries from this neighbour's Hello history (we "undo history");
- o otherwise, if nr is larger (modulo 2^{16}) than ne, then the sending node has decreased its Hello interval, and some Hellos were lost; the receiving node adds (nr - ne) 0 bits to the Hello history (we "fast-forward").

The receiving node then appends a 1 bit to the Hello history and sets ne to (nr + 1). If the Interval field of the received Hello is not zero, it resets the neighbour's hello timer to 1.5 times the advertised Interval (the extra margin allows for delay due to jitter).

Whenever either Hello timer associated to a neighbour expires, the local node adds a 0 bit to the corresponding Hello history, and increments the expected Hello number. If both Hello histories are empty (they contain 0 bits only), the neighbour entry is flushed; otherwise, the relevant hello timer is reset to the value advertised in the last Hello of that kind received from this neighbour (no extra margin is necessary in this case, since jitter was already taken into account when computing the timeout that has just expired).

A.2. Cost Computation

This section describes two algorithms suitable for computing costs (Section 3.4.3) based on Hello history. Appendix A.2.1 applies to wired links, and Appendix A.2.2 to wireless links. RECOMMENDED default values of the parameters that appear in these algorithms are given in Appendix B.

A.2.1. k-out-of-j

K-out-of-j link sensing is suitable for wired links that are either up, in which case they only occasionally drop a packet, or down, in which case they drop all packets.

The k-out-of-j strategy is parameterised by two small integers k and j, such that $0 < k \leq j$, and the nominal link cost, a constant $C \geq 1$. A node keeps a history of the last j hellos; if k or more of

those have been correctly received, the link is assumed to be up, and the rxcost is set to C; otherwise, the link is assumed to be down, and the rxcost is set to infinity.

Since Babel supports two kinds of Hellos, a Babel node performs k-out-of-j twice for each neighbour, once on the Unicast and once on the Multicast Hello history. If either of the instances of k-out-of-j indicates that the link is up, then the link is assumed to be up, and the rxcost is set to C; if both instances indicate that the link is down, then the link is assumed to be down, and the rxcost is set to infinity. In other words, the resulting rxcost is the minimum of the rxcosts yielded by the two instances of k-out-of-j link sensing.

The cost of a link performing k-out-of-j link sensing is defined as follows:

- o cost = FFFF hexadecimal if rxcost = FFFF hexadecimal;
- o cost = txcost otherwise.

A.2.2. ETX

Unlike wired links which are bimodal (either up or down), wireless links exhibit continuous variation of the link quality. Naive application of hop-count routing in networks that use wireless links for transit tends to select long, lossy links in preference to shorter, lossless links, which can dramatically reduce throughput. For that reason, a routing protocol designed to support wireless links must perform some form of link-quality estimation.

The Expected Transmission Cost algorithm, or ETX [ETX], is a simple link-quality estimation algorithm that is designed to work well with the IEEE 802.11 MAC [IEEE802.11]. By default, the IEEE 802.11 MAC performs Automatic Repeat Query (ARQ) and rate adaptation on unicast frames, but not on multicast frames, which are sent at a fixed rate with no ARQ; therefore, measuring the loss rate of multicast frames yields a useful estimate of a link's quality.

A node performing ETX link quality estimation uses a neighbour's Multicast Hello history to compute an estimate, written beta, of the probability that a Hello TLV is successfully received. Beta can be computed as the fraction of 1 bits within a small number (say, 6) of the most recent entries in the Multicast Hello history, or it can be an exponential average, or some combination of both approaches. Let rxcost be $256 / \text{beta}$.

Let α be $\text{MIN}(1, 256/\text{txcost})$, an estimate of the probability of successfully sending a Hello TLV. The cost is then computed by

$$\text{cost} = 256/(\alpha * \beta)$$

or, equivalently,

$$\text{cost} = (\text{MAX}(\text{txcost}, 256) * \text{rxcost}) / 256.$$

Since the IEEE 802.11 MAC performs ARQ on unicast frames, unicast frames do not provide a useful measure of link quality, and therefore ETX ignores the Unicast Hello history. Thus, a node performing ETX link-quality estimation will not route through neighbouring nodes unless they send periodic Multicast Hellos (possibly in addition to Unicast Hellos).

A.3. Route selection and hysteresis

Route selection (Section 3.6) is the process by which a node selects a single route among the routes that it has available towards a given destination. With Babel, any route selection procedure that only ever chooses feasible routes with a finite metric will yield a set of loop-free routes; however, in the presence of continuously variable metrics such as ETX (Appendix A.2.2), a naive route selection procedure might lead to persistent oscillations. Such oscillations can be limited or avoided altogether by implementing hysteresis within the route selection algorithm, i.e., by making the route selection algorithm sensitive to a route's history. Any reasonable hysteresis algorithm should yield good results; the following algorithm is simple to implement and has been successfully deployed in a variety of environments.

For every route R , in addition to the route's metric $m(R)$, maintain a smoothed version of $m(R)$ written $ms(R)$ (we RECOMMEND computing $ms(R)$ as an exponentially smoothed average (see Section 3.7 of [RFC793]) of $m(R)$ with a time constant equal to the Hello interval multiplied by a small number such as 3). If no route to a given destination is selected, then select the route with the smallest metric, ignoring the smoothed metric. If a route R is selected, then switch to a route R' only when both $m(R') < m(R)$ and $ms(R') < ms(R)$.

Intuitively, the smoothed metric is a long-term estimate of the quality of a route. The algorithm above works by only switching routes when both the instantaneous and the long-term estimate of the route's quality indicate that switching is profitable.

Appendix B. Protocol parameters

The choice of time constants is a trade-off between fast detection of mobility events and protocol overhead. Two instances of Babel running with different time constants will interoperate, although the resulting worst-case convergence time will be dictated by the slower of the two.

The Hello interval is the most important time constant: an outage or a mobility event is detected within 1.5 to 3.5 Hello intervals. Due to Babel's use of a redundant route table, and due to its reliance on triggered updates and explicit requests, the Update interval has little influence on the time needed to reconverge after an outage: in practice, it only has a significant effect on the time needed to acquire new routes after a mobility event. While the protocol allows intervals as low as 10ms, such low values would cause significant amounts of protocol traffic for little practical benefit.

The following values have been found to work well in a variety of environments, and are therefore RECOMMENDED default values:

Multicast Hello Interval: 4 seconds.

Unicast Hello Interval: infinite (no Unicast Hellos are sent).

Link cost: estimated using ETX on wireless links; 2-out-of-3 with C=96 on wired links.

IHU Interval: the advertised IHU interval is always 3 times the Multicast Hello interval. IHUs are actually sent with each Hello on lossy links (as determined from the Hello history), but only with every third Multicast Hello on lossless links.

Update Interval: 4 times the Multicast Hello interval.

IHU Hold Time: 3.5 times the advertised IHU interval.

Route Expiry Time: 3.5 times the advertised update interval.

Request timeout: initially 2 seconds, doubled every time a request is resent, up to a maximum of three times.

Urgent timeout: 0.2 seconds.

Source GC time: 3 minutes.

Appendix C. Route filtering

Route filtering is a procedure where an instance of a routing protocol either discards some of the routes announced by its neighbours, or learns them with a metric that is higher than what would be expected. Like all distance-vector protocols, Babel has the ability to apply arbitrary filtering to the routes it learns, and implementations of Babel that apply different sets of filtering rules will interoperate without causing routing loops. The protocol's ability to perform route filtering is a consequence of the latitude given in Section 3.5.2: Babel can use any metric that is strictly monotonic, including one that assigns an infinite metric to a selected subset of routes. (See also Section 3.8.1, where requests for nonexistent routes are treated in the same way as requests for routes with infinite metric.)

It is not in general correct to learn a route with a metric smaller than the one it was announced with, or to replace a route's destination prefix with a more specific (longer) one. Doing either of these may cause persistent routing loops.

Route filtering is a useful tool, since it allows fine-grained tuning of the routing decisions made by the routing protocol. Accordingly, some implementations of Babel implement a rich configuration language that allows applying filtering to sets of routes defined, for example, by incoming interface and destination prefix.

In order to limit the consequences of misconfiguration, Babel implementations provide a reasonable set of default filtering rules even when they don't allow configuration of filtering by the user. At a minimum, they discard routes with a destination prefix in fe80::/64, ff00::/8, 127.0.0.1/32, 0.0.0.0/32 and 224.0.0.0/8.

Appendix D. Considerations for protocol extensions

Babel is an extensible protocol, and this document defines a number of mechanisms that can be used to extend the protocol in a backwards compatible manner:

- o increasing the version number in the packet header;
- o defining new TLVs;
- o defining new sub-TLVs (with or without the mandatory bit set);
- o defining new AEs;
- o using the packet trailer.

This appendix is intended to guide designers of protocol extensions in choosing a particular encoding.

The version number in the Babel header should only be increased if the new version is not backwards compatible with the original protocol.

In many cases, an extension could be implemented either by defining a new TLV, or by adding a new sub-TLV to an existing TLV. For example, an extension whose purpose is to attach additional data to route updates can be implemented either by creating a new "enriched" Update TLV, by adding a non-mandatory sub-TLV to the Update TLV, or by adding a mandatory sub-TLV.

The various encodings are treated differently by implementations that do not understand the extension. In the case of a new TLV or of a sub-TLV with the mandatory bit set, the whole TLV is ignored by implementations that do not implement the extension, while in the case of a non-mandatory sub-TLV, the TLV is parsed and acted upon, and only the unknown sub-TLV is silently ignored. Therefore, a non-mandatory sub-TLV should be used by extensions that extend the Update in a compatible manner (the extension data may be silently ignored), while a mandatory sub-TLV or a new TLV must be used by extensions that make incompatible extensions to the meaning of the TLV (the whole TLV must be thrown away if the extension data is not understood).

Experience shows that the need for additional data tends to crop up in the most unexpected places. Hence, it is recommended that extensions that define new TLVs should make them self-terminating, and allow attaching sub-TLVs to them.

Adding a new AE is essentially equivalent to adding a new TLV: Update TLVs with an unknown AE are ignored, just like unknown TLVs. However, adding a new AE is more involved than adding a new TLV, since it creates a new set of compression state. Additionally, since the Next Hop TLV creates state specific to a given address family, as opposed to a given AE, a new AE for a previously defined address family must not be used in the Next Hop TLV if backwards compatibility is required. A similar issue arises with Update TLVs with unknown AEs establishing a new router-id (due to the Router-Id flag being set). Therefore, defining new AEs must be done with care if compatibility with unextended implementations is required.

The packet trailer is intended to carry cryptographic signatures that only cover the packet body; storing the cryptographic signatures in the packet trailer avoids clearing the signature before computing a hash of the packet body, and makes it possible to check a

cryptographic signature before running the full, stateful TLV parser. Hence, only TLVs that don't need to be protected by cryptographic security protocols should be allowed in the packet trailer. Any such TLVs should be easy to parse, and in particular should not require stateful parsing.

Appendix E. Stub Implementations

Babel is a fairly economic protocol. Updates take between 12 and 40 octets per destination, depending on the address family and how successful compression is; in a double-stack flat network, an average of less than 24 octets per update is typical. The route table occupies about 35 octets per IPv6 entry. To put these values into perspective, a single full-size Ethernet frame can carry some 65 route updates, and a megabyte of memory can contain a 20000-entry route table and the associated source table.

Babel is also a reasonably simple protocol. One complete implementation consists of less than 12 000 lines of C code, and it compiles to less than 120 kB of text on a 32-bit CISC architecture; about half of this figure is due to protocol extensions and user-interface code.

Nonetheless, in some very constrained environments, such as PDAs, microwave ovens, or abacuses, it may be desirable to have subset implementations of the protocol.

There are many different definitions of a stub router, but for the needs of this section a stub implementation of Babel is one that announces one or more directly attached prefixes into a Babel network but doesn't reannounce any routes that it has learnt from its neighbours, and always prefers the direct route to a directly attached prefix to a route learned over the Babel protocol, even when the prefixes are the same. It may either maintain a full routing table, or simply select a default gateway through any one of its neighbours that announces a default route. Since a stub implementation never forwards packets except from or to a directly attached link, it cannot possibly participate in a routing loop, and hence it need not evaluate the feasibility condition or maintain a source table.

No matter how primitive, a stub implementation must parse sub-TLVs attached to any TLVs that it understands and check the mandatory bit. It must answer acknowledgment requests and must participate in the Hello/IHU protocol. It must also be able to reply to seqno requests for routes that it announces and, and it should be able to reply to route requests.

Experience shows that an IPv6-only stub implementation of Babel can be written in less than 1000 lines of C code and compile to 13 kB of text on 32-bit CISC architecture.

Appendix F. Compatibility with previous versions

The protocol defined in this document is a successor to the protocol defined in [RFC6126] and [RFC7557]. While the two protocols are not entirely compatible, the new protocol has been designed so that it can be deployed in existing RFC 6126 networks without requiring a flag day.

There are three optional features that make this protocol incompatible with its predecessor. First of all, RFC 6126 did not define Unicast hellos (Section 3.4.1), and an implementation of RFC 6126 will mis-interpret a Unicast Hello for a Multicast one; since the sequence number space of Unicast Hellos is distinct from the sequence space of Multicast Hellos, sending a Unicast Hello to an implementation of RFC 6126 will confuse its link quality estimator. Second, RFC 6126 did not define unscheduled Hellos, and an implementation of RFC 6126 will mis-parse Hellos with an interval equal to 0. Finally, RFC 7557 did not define mandatory sub-TLVs (Section 4.4), and thus, an implementation of RFCs 6126 and 7557 will not correctly ignore a TLV that carries an unknown mandatory sub-TLV; depending on the sub-TLV, this might cause routing pathologies.

An implementation of this specification that never sends Unicast or unscheduled Hellos and doesn't implement any extensions that use mandatory sub-TLVs is safe to deploy in a network in which some nodes implement the protocol described in RFCs 6126 and 7557.

Two changes need to be made to an implementation of RFCs 6126 and 7557 so that it can safely interoperate in all cases with implementations of this protocol. First, it needs to be modified to either ignore or process Unicast and unscheduled Hellos. Second, it needs to be modified to parse sub-TLVs of all the TLVs that it understands and that allow sub-TLVs, and to ignore the TLV if an unknown mandatory sub-TLV is found. It is not necessary to parse unknown TLVs, as these are ignored in any case.

There are other changes, but these are not of a nature to prevent interoperability:

- o the conditions on route acquisition (Section 3.5.3) have been relaxed;
- o route selection should no longer use the route's sequence number (Section 3.6);

- o the format of the packet trailer has been defined (Section 4.2);
- o router-ids with a value of all-zeros or all-ones have been forbidden (Section 4.1.3);
- o the compression state is now specific to an address family rather than an address encoding (Section 4.5);
- o packet pacing is now recommended (Section 3.1).

Appendix G. Changes from previous versions

[RFC Editor: Please delete this section before publication.]

G.1. Changes since RFC 6126

- o Changed UDP port number to 6696.
- o Consistently use router-id rather than id.
- o Clarified that the source garbage collection timer is reset after sending an update even if the entry was not modified.
- o In section "Seqno Requests", fixed an erroneous "route request".
- o In the description of the Seqno Request TLV, added the description of the Router-Id field.
- o Made router-ids all-0 and all-1 forbidden.

G.2. Changes since draft-ietf-babel-rfc6126bis-00

- o Added security considerations.

G.3. Changes since draft-ietf-babel-rfc6126bis-01

- o Integrated the format of sub-TLVs.
- o Mentioned for each TLV whether it supports sub-TLVs.
- o Added Appendix D.
- o Added a mandatory bit in sub-TLVs.
- o Changed compression state to be per-AF rather than per-AE.
- o Added implementation hint for the routing table.

- o Clarified how router-ids are computed when bit 0x40 is set in Updates.
- o Relaxed the conditions for sending requests, and tightened the conditions for forwarding requests.
- o Clarified that neighbours should be acquired at some point, but it doesn't matter when.

G.4. Changes since draft-ietf-babel-rfc6126bis-02

- o Added Unicast Hellos.
- o Added unscheduled (interval-less) Hellos.
- o Changed Appendix A to consider Unicast and unscheduled Hellos.
- o Changed Appendix B to agree with the reference implementation.
- o Added optional algorithm to avoid the hold time.
- o Changed the table of pending seqno requests to be indexed by router-id in addition to prefixes.
- o Relaxed the route acquisition algorithm.
- o Replaced minimal implementations by stub implementations.
- o Added acknowledgments section.

G.5. Changes since draft-ietf-babel-rfc6126bis-03

- o Clarified that all the data structures are conceptual.
- o Made sending and receiving Multicast Hellos a SHOULD, avoids expressing any opinion about Unicast Hellos.
- o Removed opinion about Multicast vs. Unicast Hellos (Appendix A.4).
- o Made hold-time into a SHOULD rather than MUST.
- o Clarified that Seqno Requests are for a finite-metric Update.
- o Clarified that sub-TLVs Pad1 and PadN are allowed within any TLV that allows sub-TLVs.
- o Updated IANA Considerations.

- o Updated Security Considerations.
- o Renamed routing table back to route table.
- o Made buffering outgoing updates a SHOULD.
- o Weakened advice to use modified EUI-64 in router-ids.
- o Added information about sending requests to Appendix B.
- o A number of minor wording changes and clarifications.

G.6. Changes since draft-ietf-babel-rfc6126bis-03

Minor editorial changes.

G.7. Changes since draft-ietf-babel-rfc6126bis-04

- o Renamed isotonicity to left-distributivity.
- o Minor clarifications to unicast hellos.
- o Updated requirements boilerplate to RFC 8174.
- o Minor editorial changes.

G.8. Changes since draft-ietf-babel-rfc6126bis-05

- o Added information about the packet trailer, now that it is used by draft-ietf-babel-hmac.

G.9. Changes since draft-ietf-babel-rfc6126bis-06

- o Added references to security documents.

G.10. Changes since draft-ietf-babel-rfc6126bis-07

- o Added list of obsoleted drafts to the abstract.
- o Updated references.

G.11. Changes since draft-ietf-babel-rfc6126bis-08

- o Added recommendation that route selection should not take seqnos into account.

G.12. Changes since draft-ietf-babel-rfc6126bis-09

- o Editorial changes only.

G.13. Changes since draft-ietf-babel-rfc6126bis-10

- o Editorial changes only.

G.14. Changes since draft-ietf-babel-rfc6126bis-11

- o Added recommendation that control traffic should be carried over IPv6 only.

G.15. Changes since draft-ietf-babel-rfc6126bis-12

- o Removed appendix about software availability.
- o Expanded appendix about recommended values and added more references to it in the body of the document.
- o Added appendix about route filtering.
- o Clarified definition of mandatory bit.
- o Added recommendations for packet pacing.
- o Made time limiting of full updates a SHOULD.
- o Normative language in a few more places.
- o Removed normative language from stub implementations.
- o Added requirement to clear the undefined bits in an Update.
- o Added error checking requirements.
- o Reworked security considerations.
- o Added "in octets" and "in bits" in random places.
- o Inserted full IANA registries.
- o Editorial changes.

G.16. Changes since draft-ietf-babel-rfc6126bis-13

- o Added a section about compatibility with 6126.
- o Added AE registry to IANA considerations.
- o Replaced Babel-HMAC with Babel-MAC, consistent with the change in draft-ietf-babel-hmac.
- o Removed section about external sources of willingness; filtering is a better approach.
- o Added recommendation to use a cost of 96 on wired links.
- o Editorial changes.

G.17. Changes since draft-ietf-babel-rfc6126bis-14

- o Added unscheduled Hellos to compatibility considerations.
- o Created new appendix about route selection.
- o Reworked security considerations.
- o Added some comments about packet pacing and low update intervals.

G.18. Changes since draft-ietf-babel-rfc6126bis-15

- o Implementing Babel-MAC is now recommended.

G.19. Changes since draft-ietf-babel-rfc6126bis-16

- o Make the values in Appendix B normatively recommended defaults.

G.20. Changes since draft-ietf-babel-rfc6126bis-17

- o Hysteresis in route selection is now RECOMMENDED.
- o Additive metric algebra is now RECOMMENDED default.
- o 2-out-of-3 cost computation is now RECOMMENDED on LANs.
- o Reference to RFC 793 Section 3.7 added as exponential smoothing example.

G.21. Changes since draft-ietf-babel-rfc6126bis-18

- o Reserved Address Encodings 224-254 for Experimental Use, and 255 for future expansion.

G.22. Changes since draft-ietf-babel-rfc6126bis-19

- o Mention that multi-octet fields are in big-endian.
- o Minor typos and clarifications.

Authors' Addresses

Juliusz Chroboczek
IRIF, University of Paris-Diderot
Case 7014
75205 Paris Cedex 13
France

Email: jch@irif.fr

David Schinazi
Google LLC
1600 Amphitheatre Parkway
Mountain View, California 94043
USA

Email: dschinazi.ietf@gmail.com