

Network Working Group
Internet-Draft
Updates: 6126bis
(if approved)
Intended status: Standards Track
Expires: December 17, 2017

M. Boutier
J. Chroboczek
IRIF, University of Paris-Diderot
June 15, 2017

Source-Specific Routing in Babel
draft-boutier-babel-source-specific-02

Abstract

This document describes an extension to the Babel routing protocol to support source-specific routing.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on December 17, 2017.

Copyright Notice

Copyright (c) 2017 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. TODOs	3
2. Introduction and background	3
3. Data Structures	3
3.1. The Source Table	3
3.2. The Route Table	3
3.3. The Table of Pending Requests	4
4. Data Forwarding	4
5. Protocol Operation	5
5.1. Source-specific messages	5
5.2. Route Acquisition	5
5.3. Wildcard requests	6
6. Backwards compatibility	7
6.1. Loop-avoidance	8
6.2. Starvation and Blackholes	8
7. Protocol Encoding	9
7.1. Source Prefix sub-TLV	9
7.2. Source-specific Update	9
7.3. Source-specific (Route) Request	10
7.4. Source-Specific Seqno Request	10
8. IANA Considerations	10
9. Security considerations	10
10. References	10
10.1. Normative References	10
10.2. Informative References	11
Authors' Addresses	11

1. TODOs

- o Source Prefix sub-TLV type: TBD
- o check references (Section) for BABEL in 6126bis

2. Introduction and background

Source-specific routing (other known as Source Address Dependant Routing, SAD Routing or SADR) is an extension to traditional next-hop routing where packets are routed according to both their destination and their source address. This document describes the source-specific routing extension to the Babel routing protocol as defined in 6126bis [BABEL]. It notably requires the sub-TLV mandatory bit.

Background information about source-specific routing is provided in [SS-ROUTING].

3. Data Structures

This extension adds some data to the data structures maintained by a Babel node.

3.1. The Source Table

Every Babel node maintains a source table, as described in [BABEL], Section 3.2.5. A source-specific Babel node extends this table with the following field:

- o the source prefix (sprefix, splen) specifying the source address of packets to which this entry applies.

If a source table entry has a zero length source prefix (splen equals to 0), then the entry is a non-source-specific entry, and is treated just like a source table entry defined by the original Babel protocol.

With this extension the route entry contains a source which itself contains a source prefix. Notwithstanding the accidental similarity in their names, these are two very different concepts, and should not be confused.

3.2. The Route Table

Every Babel node maintains a route table, as described in [BABEL], Section 3.2.6. With this extension, this table is indexed by the 5-tuple (prefix, plen, source prefix, source plen, router-id)

obtained from the associated source table entry.

If a route table entry has a zero length source prefix, then the entry is a non-source-specific entry, and is treated just like a route table entry defined by the original Babel protocol.

3.3. The Table of Pending Requests

Every Babel node maintains a table of pending requests, as described in [BABEL], Section 3.2.7. A source-specific Babel node extends this table with the following entry:

- o the source prefix being requested.

4. Data Forwarding

In next-hop routing, if two routing table entries overlap, then one is necessarily more specific than the other; the "longest prefix rule" specifies that the most specific applicable routing table entry is chosen.

With source-specific routing, there might no longer be a most specific applicable prefix: two routing table entries might match a given packet without one necessarily being more specific than the other. Consider for example the following fragment of a routing table:

```
(2001:DB8:0:1::/64, ::/0, A)
(::/0, 2001:DB8:0:2::/64, B)
```

This specifies that all packets with destination in 2001:DB8:0:1::/64 are to be routed through A, while packets with a source in 2001:DB8:0:2::/64 are to be routed through B. A packet with source 2001:DB8:0:2::42 and destination 2001:DB8:0:1::57 matches both rules, although neither is more specific than the other. A choice is necessary, and unless the choice being made is the same on all routers in a routing domain, persistent routing loops may occur.

A Babel implementation MUST choose routing table entries by using the so-called destination-first ordering, where a routing table entry R1 is preferred to a routing table entry R2 when either R1's destination prefix is more specific than R2's, or the destination prefixes are equal and R1's source prefix is more specific than R2's. (In more formal terms, routing table entries are compared using the lexicographic product of the destination prefix ordering by the source prefix ordering.)

In practice, this means that a source-specific Babel implementation must take care that any lower layer that performs packet forwarding obey this semantics. In particular:

- o If the lower layers implement the destination-first ordering, then the Babel implementation MAY use them directly;
- o If the lower layers can hold source-specific routes, but not with the right semantics, then the Babel implementation MUST disambiguate the routing table by using a suitable disambiguation algorithm (see [SS-ROUTING] for such an algorithm);
- o If the lower layers cannot hold source-specific routes, then a Babel implementation MUST silently ignore any source-specific routes.

5. Protocol Operation

This extension does not fundamentally change the operation of the Babel protocol. We only described the fundamental differences between the original protocol and the extension in this section. The other mechanisms described in [BABEL] (Section 3) may be inferred by using pairs of (destination, source) prefixes instead of just (destination) prefixes.

5.1. Source-specific messages

A route of this extension with a zero-length source prefix is the same than a route without source prefix (a route of the classical Babel). In both of the cases, packets are accepted independantly of their source address. Thus, a route is said source-specific only if its source prefix has a non-zero length.

Three messages are used to communicate informations on routes: Updates, Route Requests and Seqno Requests. With this extension, these messages carry an additionnal source prefix if (and only if) the corresponding route is source-specific. More formally, an Update, a Route Request and a Seqno Request MUST carry a source prefix if they concern a source-specific route (non-zero length source prefix) and MUST NOT carry a source prefix otherwise (zero length source prefix). A message which carry a source prefix is said source-specific.

5.2. Route Acquisition

When a non-source-specific Babel node receives a source-specific update, it just ignores it.

On receipt of a source-specific update (id, prefix, source prefix,

seqno, metric), a source-specific Babel node behaves as described in [BABEL] Section 3.5.4 though indexing entries by (neigh, id, prefix, source prefix). When a source-specific Babel node receives a non-source-specific update, it MUST consider this update as carrying a zero length source prefix.

5.3. Wildcard requests

TODO: behaviour to be defined.

5.3.1. Proposal 1

The original Babel protocol states that when a node receives a wildcard route request, it SHOULD send a full routing table dump. This extension does not change this statement: a source-specific node SHOULD send a full routing table dump when receiving a wildcard request.

Source-specific wildcard requests does not exist: a wildcard request SHOULD NOT carry a source prefix.

5.3.2. Proposal 2

We assume that a mandatory sub-TLV has a corresponding non-mandatory sub-TLV. This proposal is like Proposal 3 but instead of having multiple wildcard request TLVs, one for each kind of routes understood, we use one wildcard request with sub-TLVs corresponding to the extension. To have a full routing table dump, a node sends a wildcard requests with a non-mandatory Source sub-TLV.

A source-specific node SHOULD always attach a non-mandatory Source sub-TLV to its wildcard requests.

This proposal has been rejected because it implies to share the space of non-mandatory and mandatory sub-TLVs.

5.3.3. Proposal 3

The Babel protocol provides the ability to request a full routing table dump by sending a "wildcard request", a route request with the AE field set to 0. As the original protocol has no source-specific routes, such a request may only concern non-source-specific routes. This extension does not modify the semantics of wildcard requests in that sense: a wildcard request prompts the receiver to send its non-source-specific routes only, and a Babel node SHOULD NOT send any source-specific updates in reply to a wildcard request.

To obtain a dump of the source-specific routes, a source-specific

wildcard request **MUST** be used. A source-specific wildcard request is a wildcard request carrying a zero length source prefix.

When a node receives a source-specific wildcard request, it **SHOULD** send a dump of its routes which are source-specific "only". It **SHOULD NOT** send any non-source-specific routes in reply to a source-specific wildcard request. It **SHOULD NOT** send any source-specific routes which are under the effect of a future extension. Such extension should detail how to handle the possible combinations.

In consequence, a node requiring a full routing table dump must send both a non-source-specific wildcard request and a source-specific wildcard request.

5.3.4. Proposal 4

Wildcard requests are deprecated. Either deprecate it in 6126bis, or say the following.

A node receiving a wildcard request **SHOULD** ignore it.

This proposal has been rejected because wildcard requests speeds up the convergence of the network on boot. This is considered important.

5.3.5. Note on Overhead between (1) and (3)

Sending one wildcard request (1) instead of a few something-specific wildcard requests (3) in a negligible gain.

Non-source-specific nodes sending requests to source-specific nodes may reduce the global overhead with (3). But, if the network has no source-specific route, there is no overhead to reduce; if there is only a few source-specific routes (like in a home network), the overhead would be negligible. Thus, the interesting case is when there is a lot of source-specific routes.

We can imagine a network with a source-specific backbone announcing a default route and catching all traffic. Good old routers not supporting this extensions would be put at some backbone leafs. Is sbabeld part of that use case ?

Couldn't we just send a Route Request for **default** ?

6. Backwards compatibility

The protocol extension defined in this document is, to a great

extent, interoperable with the base protocol defined in [BABEL] (and all its known extensions). More precisely, if non-source-specific routers and source-specific routers are mixed in a single routing domain, Babel's loop-avoidance properties are preserved, and, in particular, no persistent routing loops will occur.

TODO: Should we put a warning to say it's not the case with the Experimental Track Babel ?

6.1. Loop-avoidance

The extension defined in this protocol uses a new Mandatory sub-TLV to carry the source prefix information. As discussed in Section 4.4 of [BABEL], this encoding ensures that non-source-specific routers will silently ignore the whole TLV, which is necessary to avoid persistent routing loops in hybrid networks.

Consider two nodes A and B, with A source-specific announcing a route to (D, S). Suppose that B ignores the source prefix information when it receives the update, and reannounces it as D. This is reannounced to A, which treats it as (D, ::/0). Packets destined to D but not sourced in S will be forwarded by A to B, and by B to A, causing a persistent routing loop:

```

      (D,S)                (D)
      <--                  <--
----- A ----- B
      -->
      (D,::/0)

```

6.2. Starvation and Blackholes

In general, discarding of source-specific routes by non-source-specific routers will cause routing starvation. Intuitively, unless there are enough non-source-specific routes in the network, non-source-specific routers will suffer starvation, and discard packets for destinations that are only announced by source-specific routers.

A simple yet sufficient condition for avoiding starvation is to build a connected source-specific backbone that includes all of the edge routers, and announce a (non-source-specific) default route towards the backbone. However, introducing such a default route in the network may in the same time introduce a blackhole. This tradeoff is let to the administrator.

7. Protocol Encoding

This extension defines a new sub-TLV used to carry a source prefix by the three following existing messages: Update, Route Request and Seqno Request.

7.1. Source Prefix sub-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 = TBD | Length | Source Plen | Source Prefix...
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Fields:

Type Set to TBD to indicate a Source Prefix sub-TLV.

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

Source Plen The length of the advertised source prefix. This MUST NOT be 0.

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

The Source Prefix field's encoding is the same than the Prefix's. It is defined by the AE field of the corresponding TLV.

Remark that this sub-TLV is a Mandatory sub-TLV. The whole TLV MUST be ignored if that TLV is not recognized. Otherwise, routing loops may occur.

7.2. Source-specific Update

The source-specific Update is an Update TLV with a Source Prefix sub-TLV. It advertises or retracts source-specific routes in the same manner than routes with non-source-specific Updates (see [BABEL]) except for wildcard updates.

Wildcard updates MUST NOT carry any source prefix. Wildcard updates (in fact, wildcard retraction) are used when a Babel node stops: a receiver retracts all routes announced by the announcing node. There is no use case for source-specific wildcard updates. A source-specific Babel node receiving a (legacy) wildcard update MUST retract all routes it learns from this node (including source-specific ones).

Contrary to the destination prefix, this extension does not compress the source prefix attached to Updates. The destination prefix uses compression as defined in [BABEL] for Updates with Mandatory

extensions.

However, as defined in [BABEL] (Section 4.5), the compression is allowed for the destination prefix of source-specific routes. Legacy implementation will correctly update their parser state, while ignoring the whole TLV afterwards.

7.3. Source-specific (Route) Request

TODO: A source-specific Route Request prompts the receiver to send an update for a given pair of destination and source prefixes. It MUST NOT be used to request a full routing table dump. The Source Prefix sub-TLV of a wildcard source-specific Route Request (Request with AE equals to 0 and a Source Prefix sub-TLV) MIGHT be ignored: a receiver MIGHT reply by a full routing table dump.

7.4. Source-Specific Seqno Request

A source-specific Seqno Request is just like a Seqno Request for a source-specific route. It uses the same mechanisms described in [BABEL].

8. IANA Considerations

IANA is instructed to add the following entry to the "Babel sub-TLV Types" registry:

+-----+-----+-----+		
Type	Name	Reference
+-----+-----+-----+		
TBD	Source Prefix	(this document)
+-----+-----+-----+		

9. Security considerations

The extension defined in this document adds a new sub-TLV to three TLVs already present in the original Babel protocol. It does not by itself change the security properties of the protocol.

10. References

10.1. Normative References

[BABEL] Chroboczek, J., "The Babel Routing Protocol", Internet Draft draft-ietf-babel-rfc6126bis-02, May 2017.

10.2. Informative References

[SS-ROUTING]

Boutier, M. and J. Chroboczek, "Source-Specific Routing",
August 2014.

In Proc. IFIP Networking 2015. A slightly earlier
version is available online from
<http://arxiv.org/pdf/1403.0445>.

Authors' Addresses

Matthieu Boutier
IRIF, University of Paris-Diderot
Case 7014
75205 Paris Cedex 13,
France

Email: boutier@irif.fr

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

Email: jch@irif.fr

Network Working Group
Internet-Draft

Updates: 6126bis (if approved)
Intended status: Standards Track
Expires: January 4, 2018

M. Boutier
J. Chroboczek
IRIF, University of Paris-Diderot
July 03, 2017

Source-Specific Routing in Babel
draft-boutier-babel-source-specific-03

Abstract

This document describes an extension to the Babel routing protocol to support source-specific routing.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on January 4, 2018.

Copyright Notice

Copyright (c) 2017 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. TODOs	2
2. Introduction and background	2
3. Data Structures	3
3.1. The Source Table	3
3.2. The Route Table	3
3.3. The Table of Pending Requests	3
4. Data Forwarding	3
5. Protocol Operation	5
5.1. Source-specific messages	5
5.2. Route Acquisition	5
5.3. Wildcard retractions (update)	5
5.4. Wildcard requests	6
6. Compatibility with the base protocol	8
6.1. Loop-avoidance	8
6.2. Starvation and Blackholes	8
7. Protocol Encoding	9
7.1. Source Prefix sub-TLV	9
7.2. Source-specific Update	9
7.3. Source-specific (Route) Request	10
7.4. Source-Specific Seqno Request	10
8. IANA Considerations	10
9. Security considerations	10
10. References	10
10.1. Normative References	10
10.2. Informative References	11
Authors' Addresses	11

1. TODOs

- o Source Prefix sub-TLV type: TBD
- o check references (Section) for BABEL in 6126bis
- o define wildcard Requests behaviour

2. Introduction and background

Source-specific routing (also known as Source Address Dependant Routing, SAD Routing or SADR) is an extension to traditional next-hop routing where packets are routed according to both their destination and their source address. This document describes the source-specific routing extension to the Babel routing protocol as defined in 6126bis [BABEL].

Background information about source-specific routing is provided in [SS-ROUTING].

3. Data Structures

This extension adds some data to the data structures maintained by a Babel node.

3.1. The Source Table

Every Babel node maintains a source table, as described in [BABEL], Section 3.2.5. A source-specific Babel node extends this table with the following field:

- o the source prefix (sprefix, splen) specifying the source address of packets to which this entry applies.

If a source table entry has a zero length source prefix (splen equals to 0), then the entry is a non-source-specific entry, and is treated just like a source table entry defined by the original Babel protocol.

With this extension the route entry contains a source which itself contains a source prefix. These are two very different concepts, and should not be confused.

3.2. The Route Table

Every Babel node maintains a route table, as described in [BABEL], Section 3.2.6. With this extension, this table is indexed by the 5-tuple (prefix, plen, source prefix, source plen, router-id) obtained from the associated source table entry.

If a route table entry has a zero length source prefix, then the entry is a non-source-specific entry, and is treated just like a route table entry defined by the original Babel protocol.

3.3. The Table of Pending Requests

Every Babel node maintains a table of pending requests, as described in [BABEL], Section 3.2.7. A source-specific Babel node extends this table with the following entry:

- o the source prefix being requested.

4. Data Forwarding

In next-hop routing, if two routing table entries overlap, then one is necessarily more specific than the other; the "longest prefix rule" specifies that the most specific applicable routing table entry is chosen.

With source-specific routing, there might no longer be a most specific applicable prefix: two routing table entries might match a given packet without one necessarily being more specific than the other. Consider for example the following fragment of a routing table:

```
(2001:DB8:0:1::/64, ::/0, A)
```

```
(::/0, 2001:DB8:0:2::/64, B)
```

This specifies that all packets with destination in 2001:DB8:0:1::/64 are to be routed through A, while packets with a source in 2001:DB8:0:2::/64 are to be routed through B. A packet with source 2001:DB8:0:2::42 and destination 2001:DB8:0:1::57 matches both rules, although neither is more specific than the other. A choice is necessary, and unless the choice being made is the same on all routers in a routing domain, persistent routing loops may occur. More informations are available in [SS-ROUTING] Section IV.C.

A Babel implementation **MUST** choose routing table entries by using the so-called destination-first ordering, where a routing table entry R1 is preferred to a routing table entry R2 when either R1's destination prefix is more specific than R2's, or the destination prefixes are equal and R1's source prefix is more specific than R2's. (In more formal terms, routing table entries are compared using the lexicographic product of the destination prefix ordering by the source prefix ordering.)

In practice, this means that a source-specific Babel implementation must take care that any lower layer that performs packet forwarding obey this semantics. In particular:

- o If the lower layers implement the destination-first ordering, then the Babel implementation **MAY** use them directly;
- o If the lower layers can hold source-specific routes, but not with the right semantics, then the Babel implementation **MUST** disambiguate the routing table by using a suitable disambiguation algorithm (see [SS-ROUTING] Section V.B for such an algorithm);
- o If the lower layers cannot hold source-specific routes, then a Babel implementation **MUST** silently ignore (drop) any source-specific routes.

5. Protocol Operation

This extension does not fundamentally change the operation of the Babel protocol. We only describe the fundamental differences between the original protocol and the extension in this section. The other mechanisms described in [BABEL] (Section 3) are extended to pairs of (destination, source) prefixes instead of just (destination) prefixes.

5.1. Source-specific messages

Three messages are used to communicate informations on routes: Updates, Route Requests and Seqno Requests. With this extension, these messages carry an additionnal source prefix if (and only if) the corresponding route is source-specific. More formally, an Update, a Route Request and a Seqno Request MUST carry a source prefix if they concern a source-specific route (non-zero length source prefix) and MUST NOT carry a source prefix otherwise (zero length source prefix). A message which carries a source prefix is said to be source-specific.

5.2. Route Acquisition

When a non-source-specific Babel node receives a source-specific update, it silently ignores it.

TODO{On receipt of a source-specific update (id, prefix, source prefix, seqno, metric), a source-specific Babel node behaves as described in [BABEL] Section 3.5.4 though indexing entries by (neigh, id, prefix, source prefix).} When a source-specific Babel node receives a non-source-specific update, it MUST treat this update as carrying a zero length source prefix.

5.3. Wildcard retractions (update)

The original protocol defines a wildcard update with AE equals to 0 as being a wildcard retraction. A node receiving a wildcard retraction on an interface must consider that the sending node retracts all the routes it advertised on this interface.

Wildcard retractions are used when a node is about to leave the network. Thus, this extension does not define source-specific wildcard retraction, but extends wildcard retraction to apply also to source-specific routes. More formally, a wildcard update MUST NOT carry a source prefix, and a source-specific Babel node receiving a (legacy) wildcard update MUST retract all routes it learns from this node (including source-specific ones).

5.4. Wildcard requests

TODO: behaviour to be defined.

5.4.1. Proposal 1

The original Babel protocol states that when a node receives a wildcard route request, it SHOULD send a full routing table dump. This extension does not change this statement: a source-specific node SHOULD send a full routing table dump when receiving a wildcard request.

Source-specific wildcard requests does not exist: a wildcard request SHOULD NOT carry a source prefix.

5.4.2. Proposal 2

We assume that a mandatory sub-TLV has a corresponding non-mandatory sub-TLV. This proposal is like Proposal 3 but instead of having multiple wildcard request TLVs, one for each kind of route understood, we use one wildcard request with sub-TLVs corresponding to the extension. To have a full routing table dump, a node sends a wildcard requests with a non-mandatory Source sub-TLV.

A source-specific node SHOULD always attach a non-mandatory Source sub-TLV to its wildcard requests.

This proposal has been rejected because it implies to share the space of non-mandatory and mandatory sub-TLVs.

5.4.3. Proposal 3 (mentionned by Juliusz)

The Babel protocol provides the ability to request a full routing table dump by sending a "wildcard request", a route request with the AE field set to 0. As the original protocol has no source-specific routes, such a request may only concern non-source-specific routes. This extension does not modify the semantics of wildcard requests in that sense: a wildcard request prompts the receiver to send its non-source-specific routes only, and a Babel node SHOULD NOT send any source-specific updates in reply to a wildcard request.

To obtain a dump of the source-specific routes, a source-specific wildcard request MUST be used. A source-specific wildcard request is a wildcard request carrying a zero length source prefix.

When a node receives a source-specific wildcard request, it SHOULD send a dump of its routes which are source-specific "only". It SHOULD NOT send any non-source-specific routes in reply to a source-

specific wildcard request. It SHOULD NOT send any source-specific routes which are under the effect of a future extension. Such extension should detail how to handle the possible combinations.

In consequence, a node requiring a full routing table dump must send both a non-source-specific wildcard request and a source-specific wildcard request.

5.4.4. Proposal 4 (mentionned by Juliusz)

Wildcard requests are deprecated. Either deprecate it in 6126bis, or say the following.

A node receiving a wildcard request SHOULD ignore it.

This proposal has been rejected because wildcard requests speeds up the convergence of the network on boot. This is considered important.

5.4.5. Proposal 5 (mentionned by David)

By default, a vanilla wildcard request triggers a dump of all non-specific routes. We define a new non-mandatory sub-TLV on Route Requests called "Requested Route Types" that contains an array of all the types of routes this request is requesting.

```

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 = TBD | Length | RR Type 1 | RR Type 2...
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

We also create a registry of Requested Route (RR) types, for example:

0 = Regular

1 = Source-Specific

2 = TOS-specific

etc.

A node receiving a Requested Route Types sub-TLV in a wildcard request SHOULD send back a dump of all its routes corresponding to the requested types or to a combination of these types.

6. Compatibility with the base protocol

The protocol extension defined in this document is, to a great extent, interoperable with the base protocol defined in [BABEL] (and all its known extensions). More precisely, if non-source-specific routers and source-specific routers are mixed in a single routing domain, Babel's loop-avoidance properties are preserved, and, in particular, no persistent routing loops will occur.

However, this extension is not compatible with the Experimental Track's Babel Routing Protocol (RFC 6126). It requires the mandatory sub-TLV introduced in [BABEL]. Consequently, this extension **MUST NOT** be used with routers implementing RFC 6126, otherwise persistent routing loops may occur.

6.1. Loop-avoidance

The extension defined in this protocol uses a new Mandatory sub-TLV to carry the source prefix information. As discussed in Section 4.4 of [BABEL], this encoding ensures that non-source-specific routers will silently ignore the whole TLV, which is necessary to avoid persistent routing loops in hybrid networks.

Consider two nodes A and B, with A source-specific announcing a route to (D, S). Suppose that B merely ignores the source prefix information when it receives the update rather than ignoring the sub-TLV, and reannounces the route as D. This reannouncement reaches A, which treats it as (D, ::/0). Packets destined to D but not sourced in S will be forwarded by A to B, and by B to A, causing a persistent routing loop:

```

      (D,S)                (D)
      <--                  <--
----- A ----- B
      -->
      (D,::/0)

```

6.2. Starvation and Blackholes

In general, discarding source-specific routes by non-source-specific routers will cause route starvation. Intuitively, unless there are enough non-source-specific routes in the network, non-source-specific routers will suffer starvation, and discard packets for destinations that are only announced by source-specific routers.

A simple yet sufficient condition for avoiding starvation is to build a connected source-specific backbone that includes all of the edge

routers, and announce a (non-source-specific) default route towards the backbone.

7. Protocol Encoding

This extension defines a new sub-TLV used to carry a source prefix by the three following existing messages: Update, Route Request and Seqno Request.

7.1. Source Prefix sub-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 = TBD   |      Length      |  Source Plen  |  Source Prefix...
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Fields:

Type Set to TBD to indicate a Source Prefix sub-TLV.

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

Source Plen The length of the advertised source prefix. This MUST NOT be 0.

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

The Source Prefix field's encoding (AE) is the same as the Prefix's. It is defined by the AE field of the corresponding TLV.

Note that this sub-TLV is a Mandatory sub-TLV. The whole TLV MUST be ignored if that TLV is not recognized as described in Section 4.4. Otherwise, routing loops may occur.

7.2. Source-specific Update

The source-specific Update is an Update TLV with a Source Prefix sub-TLV. It advertises or retracts source-specific routes in the same manner than routes with non-source-specific Updates (see [BABEL]). This TLV MUST NOT be attached to wildcard updates.

Contrary to the destination prefix, this extension does not compress the source prefix attached to Updates. The destination prefix uses compression as defined in [BABEL] for Updates with Mandatory extensions.

However, as defined in [BABEL] (Section 4.5), the compression is allowed for the destination prefix of source-specific routes. Legacy implementation will correctly update their parser state, while ignoring the whole TLV afterwards.

7.3. Source-specific (Route) Request

TODO: A source-specific Route Request prompts the receiver to send an update for a given pair of destination and source prefixes. It MUST NOT be used to request a full routing table dump. The Source Prefix sub-TLV of a wildcard source-specific Route Request (Request with AE equals to 0 and a Source Prefix sub-TLV) MIGHT be ignored: a receiver MIGHT reply by a full routing table dump.

7.4. Source-Specific Seqno Request

A source-specific Seqno Request is just like a Seqno Request for a source-specific route. It uses the same mechanisms described in [BABEL].

8. IANA Considerations

IANA is instructed to add the following entry to the "Babel sub-TLV Types" registry:

Type	Name	Reference
TBD	Source Prefix	(this document)

9. Security considerations

The extension defined in this document adds a new sub-TLV to three TLVs already present in the original Babel protocol. It does not by itself change the security properties of the protocol.

10. References

10.1. Normative References

[BABEL] Chroboczek, J., "The Babel Routing Protocol", Internet Draft draft-ietf-babel-rfc6126bis-02, May 2017.

10.2. Informative References

[SS-ROUTING]

Boutier, M. and J. Chroboczek, "Source-Specific Routing", August 2014.

In Proc. IFIP Networking 2015. A slightly earlier version is available online from <http://arxiv.org/pdf/1403.0445>.

Authors' Addresses

Matthieu Boutier
IRIF, University of Paris-Diderot
Case 7014
75205 Paris Cedex 13
France

Email: boutier@irif.fr

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

Email: jch@irif.fr

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: November 25, 2017

J. Chroboczek
IRIF, University of Paris-Diderot
May 24, 2017

The Babel Routing Protocol
draft-ietf-babel-rfc6126bis-02

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.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on November 25, 2017.

Copyright Notice

Copyright (c) 2017 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	3
1.1. Features	3
1.2. Limitations	4
1.3. Specification of Requirements	4
2. Conceptual Description of the Protocol	4
2.1. Costs, Metrics and Neighbourship	5
2.2. The Bellman-Ford Algorithm	5
2.3. Transient Loops in Bellman-Ford	6
2.4. Feasibility Conditions	6
2.5. Solving Starvation: Sequencing Routes	8
2.6. Requests	9
2.7. Multiple Routers	10
2.8. Overlapping Prefixes	11
3. Protocol Operation	11
3.1. Message Transmission and Reception	12
3.2. Data Structures	12
3.3. Acknowledged Packets	16
3.4. Neighbour Acquisition	16
3.5. Routing Table Maintenance	18
3.6. Route Selection	22
3.7. Sending Updates	23
3.8. Explicit Route Requests	25
4. Protocol Encoding	29
4.1. Data Types	29
4.2. Packet Format	30
4.3. TLV Format	31
4.4. Sub-TLV Format	31
4.5. Parser state	32
4.6. Details of Specific TLVs	33
4.7. Details of specific sub-TLVs	43
5. IANA Considerations	43
6. Security Considerations	44
7. References	44
7.1. Normative References	44
7.2. Informative References	44
Appendix A. Cost and Metric Computation	45
A.1. Maintaining Hello History	45
A.2. Cost Computation	46
A.3. Metric Computation	47
Appendix B. Constants	48
Appendix C. Considerations for protocol extensions	49
Appendix D. Simplified Implementations	50
Appendix E. Software Availability	50
Appendix F. Changes from previous versions	51
F.1. Changes since RFC 6126	51
F.2. Changes since draft-ietf-babel-rfc6126bis-00	51

F.3. Changes since draft-ietf-babel-rfc6126bis-01	51
Author's Address	52

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.

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 prefix is originated by multiple routers, Babel may occasionally create a transient routing loop for this particular prefix; this loop disappears in a time proportional to its 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 reasonable 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 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), whichever 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, Babel does impose a hold time when a prefix is retracted (Section 3.5.5). 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. Hence, if a previously deaggregated prefix becomes aggregated, it will be unreachable for a few minutes. This makes Babel unsuitable for use in mobile 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", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC2119].

2. Conceptual Description of the Protocol

Babel is a mostly loop-free distance vector protocol: it is based on the Bellman-Ford protocol, 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

As many routing algorithms, Babel computes costs of links between any two neighbouring nodes, abstract values attached to the edges between two nodes. We write $C(A, B)$ for the cost of the edge from node A to node B .

Given a route between any two nodes, the metric of the route is the sum of the costs of all the edges along the route. The goal of the routing algorithm is to compute, for every source S , the tree of the 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 broadcasts 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.

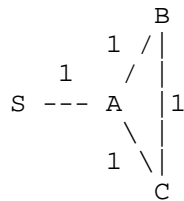
Periodically, every node B sends to all of its neighbours a route update, a message containing $D(B)$. When a neighbour A of B receives the route update, it checks whether B is its selected next hop; if that is the case, then $NH(A)$ is set to B , and $D(A)$ is set to $C(A, B) + D(B)$. If that is not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

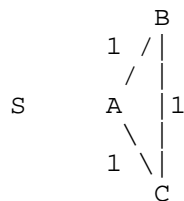
2.3. Transient Loops in Bellman-Ford

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 diagram:



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.

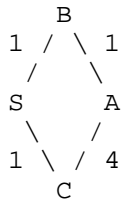
Many different feasibility conditions are possible. For example, BGP can be modelled as being a distance-vector protocol with a (rather drastic) feasibility condition: a routing update is only accepted when the receiving node's AS number is not included in the update's AS-Path attribute (note that BGP's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

Another simple feasibility condition, used in Destination-Sequenced Distance-Vector (DSDV) routing [DSDV] and in Ad hoc On-Demand Distance Vector (AODV) routing, 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 could decide that a route is feasible only 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 successor, then $D(B) < D(A)$, which implies that the forwarding graph is loop-free.

Babel uses a slightly more refined feasibility condition, used in 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 $D(A) \leq FD(A)$, $D(B) < FD(A)$.

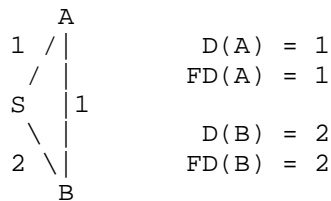
To see that it is strictly less restrictive, consider the following diagram, where A has selected the route through B, and $D(A) = FD(A) = 2$. Since $D(C) = 1 < FD(A)$, the alternate route through C is feasible for A, although its metric $C(A, C) + D(C) = 5$ is larger than that of the currently selected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



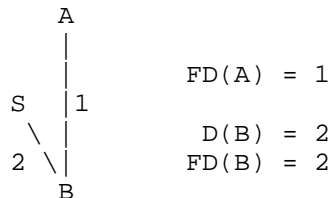
To show that this feasibility condition still guarantees loop-freedom, recall that at the time when A accepts an update from B, the metric $D(B)$ announced by B is no smaller than $FD(B)$; since it is smaller than $FD(A)$, at that point in time $FD(B) < FD(A)$. Since this property is preserved when A sends updates, 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:



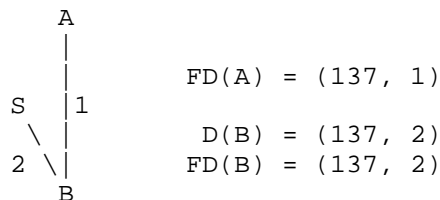
The only route available from A to S, the one that goes through B, is not feasible: A suffers from a spurious starvation.

At this point, the whole network must be rebooted in order to solve the starvation; this is essentially what EIGRP does when it performs a global synchronisation of all the routers in the network with the source (the "active" phase of EIGRP).

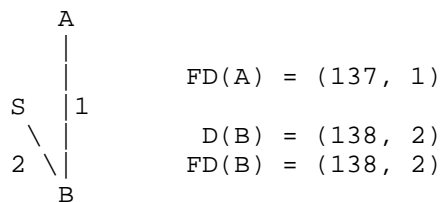
Babel reacts to starvation in a less drastic manner, by using sequenced routes, a technique introduced by DSDV and adopted by AODV. In addition to a metric, every route carries a sequence number, a nondecreasing integer that is propagated unchanged through the network and is only ever incremented by 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 necessarily used in route selection: a node will normally ignore the sequence number when selecting a route (Section 3.6).

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 duplicatedetection for the requests that they forward.

2.7. Multiple Routers

The above discussion assumes that every 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 diagram, 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.

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. Babel avoids this pitfall by maintaining an "unreachable" route for a few minutes after a route is retracted; the time for which such a route must be maintained should be the worst-case propagation time of the retraction of the route to C.

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. We suggest that router-ids should be assigned in modified EUI-64 format [ADDRARCH]. (As a matter of fact, the protocol encoding is slightly more compact when router-ids are assigned in the same manner as the IPv6 layer assigns host IDs.)

3.1. Message Transmission and Reception

Babel protocol packets are sent in the body of a UDP datagram. Each Babel packet consists of zero or more TLVs. Most TLVs may contain sub-TLVs.

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 (this is the usual case) 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 Hello TLVs and acknowledgements, all Babel TLVs can be sent to either unicast or multicast addresses, and their semantics does not depend on whether the destination was a unicast or 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 is applied to packets sent by a Babel speaker: outgoing TLVs are buffered and SHOULD be sent with a small 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 exact delay and amount of jitter applied to a packet depends on whether it contains any urgent TLVs. Acknowledgement TLVs MUST be sent before the deadline specified in the corresponding request. The particular class of updates specified in Section 3.7.2 MUST be sent in a timely manner. The particular class of request and update TLVs specified in Section 3.8.2 SHOULD be sent in a timely manner.

3.2. Data Structures

Every Babel speaker maintains a number of data structures. All of these data structures consist of familiar data types -- integers, IP addresses, etc. -- with the exception of sequence numbers.

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 serial numbers is defined as follows.

Given a seqno s and an 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 Hello seqno, a 16-bit integer that is sent with each Hello TLV on this interface and is incremented (modulo 2^{16}) whenever a Hello is sent. (Note that an interface's Hello seqno is unrelated to the node's seqno.)

There are two timers associated with each interface table entry -- the Hello timer, which governs the sending of periodic Hello and IHU packets, and the update timer, which governs the sending of periodic route updates.

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 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 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 neighbour's expected Hello sequence number, an integer modulo 2^{16} .

There are two timers associated with each neighbour entry -- the hello timer, which is initialised from the interval value carried by Hello TLVs, and the IHU timer, which is initialised to a small multiple of the interval carried in IHU TLVs.

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 fairly common situation 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, 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 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 3.5.4.

Of course, the data structure described above is conceptual: actual implementations will likely use a different data structure, for example a table of installed routes and a set of redundant ones, or some more complicated data structure.

3.2.7. The Table of Pending Requests

The table of pending 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 prefixes, and every entry in this table contains the following data:

- o the prefix, 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 request; it governs both the resending of requests and their expiry.

3.3. Acknowledged Packets

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

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

When to request acknowledgements is a matter of local policy; the simplest strategy is to never request acknowledgements and to rely on periodic updates to ensure that any reachable routes are eventually propagated throughout the routing domain. For increased efficiency, we suggest that acknowledged packets should be used in order to send urgent updates (Section 3.7.2) 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 acknowledgement requests is not necessary, and not even recommended, as the acknowledgements cause additional traffic and may force additional Address Resolution Protocol (ARP) or Neighbour Discovery 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 an implementation detail; 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 history becomes empty (see Appendix A.2).

3.4.1. Reverse Reachability Detection

Every Babel node sends periodic Hellos over each of its interfaces. Each Hello TLV carries an increasing (modulo 2^{16}) sequence number and the interval between successive periodic packets sent on this particular interface.

In addition to the periodic Hello packets, a node MAY send unscheduled Hello packets, e.g., to accelerate link cost estimation when a new neighbour is discovered, or when link conditions have suddenly changed.

A node MAY change its Hello interval. The Hello interval MAY be decreased at any time; it SHOULD NOT be increased, except immediately before sending a Hello packet. (Equivalently, a node SHOULD send an unscheduled 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. A possible 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 SHOULD be sent to a multicast address in order to avoid an ARP or Neighbour Discovery exchange and to aggregate multiple IHUs in 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.

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 updates the value of the sending neighbour's txcost (transmission cost), from its

perspective, to the value contained in the IHU, and resets this neighbour's IHU timer to a small multiple of the value received in the IHU.

When a neighbour's IHU timer expires, its 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 -- namely, 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.

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 hellos were received recently, then the cost is infinite;
- o if the txcost is infinite, then the cost is infinite.

Note that while this document does not constrain cost computation any further, not all cost computation strategies will give good results. We give a few examples of strategies for computing a link's cost that are known to work well in practice in Appendix A.2.

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 treated as a retraction, depending on some other conditions (Section 3.5.4). If the feasibility condition is satisfied, then the update cannot possibly cause a routing loop, and the update is accepted.

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 with finite metrics large enough to cause a loop are discarded.

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 (p, 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 (p, 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 (router-id, prefix, plen); 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.

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 and the metric 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 isotonic: if $m \leq m'$, then $M(c, m) \leq M(c, m')$.

Note that while strict monotonicity is essential to the integrity of the network (persistent routing loops may appear if it is not satisfied), isotonicity is not: if it is not satisfied, Babel will still converge to a locally optimal routing table, but might not reach a global optimum (in fact, such a global optimum may not even exist).

As with cost computation, not all strategies for computing route metrics will give good results. In particular, some metrics are more likely than others to lead to routing instabilities (route flapping). In Appendix A.3, we give a number of examples of strictly monotonic, isotonic routing metrics that are known to work well in practice.

3.5.3. 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. An Update TLV itself only contains the prefix, seqno, and metric, while the next hop is derived either from the network-layer source address of the packet or from an explicit Next Hop TLV

in the same packet. The router-id is derived from a separate Router-Id TLV in the same packet, which optimises the case when multiple updates are sent with the same router-id.

Additionally, a prefix of the advertised prefix can be omitted in an Update TLV, in which case it is copied from a previous Update TLV in the same packet -- this is known as address compression [PACKETBB].

Finally, as a special optimisation for the case when a router-id coincides with the interface-id part of an IPv6 address, the router-id can optionally be derived from the low-order bits of the advertised prefix.

The encoding of updates is described in detail in Section 4.6.

3.5.4. Route Acquisition

When a Babel node receives an update (router-id, prefix, seqno, metric) from a neighbour neigh with a link cost value equal to cost, it checks whether it already has a routing table entry indexed by (neigh, router-id, prefix).

If no such entry exists:

- o if the update is unfeasible, it is ignored;
- o if the metric is infinite (the update is a retraction), the update is ignored;
- o otherwise, a new route table entry is created, indexed by (neigh, router-id, prefix), with 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 installed and the update is unfeasible, then the behaviour depends on whether the router-ids of the two entries match. If the router-ids are different, the update is treated as though it were a retraction (i.e., as though the metric were FFFF hexadecimal). If the router-ids are equal, the update is ignored;
- o otherwise (i.e., if either the update is feasible or the entry is not currently installed), then the entry's sequence number, advertised metric, metric, and router-id are updated and, unless the advertised metric is infinite, the route's expiry timer is reset to a small multiple of the Interval value included in the update.

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

3.5.5. Hold Time

When a prefix *p* is retracted, because all routes are unfeasible, too old, or have an infinite metric, 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 is retracted, a routing table entry with infinite metric is maintained as described in Section 3.5.4 above, and packets destined for that prefix MUST NOT be forwarded by following a route for a shorter prefix. The infinite metric entry is maintained until it is superseded by a feasible update; if no such update arrives within the route hold time, the entry is flushed.

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 protocol'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.

Note, however, that Babel does not naturally guarantee the stability of routing, and configuring conflicting route selection policies on different routers may lead to persistent route oscillation.

Defining a good route selection policy for Babel is an open research problem. Route selection can take into account multiple mutually contradictory criteria; in roughly decreasing order of importance, these are:

- o routes with a small metric should be preferred over routes with a large metric;
- o switching router-ids should be avoided;
- o routes through stable neighbours should be preferred over routes through unstable ones;
- o stable routes should be preferred over unstable ones;
- o switching next hops should be avoided.

A simple strategy is to choose the feasible route with the smallest metric, with a small amount of hysteresis applied to avoid switching router-ids.

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 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 sends 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., the address of a local interface, the prefix of a directly attached network, or redistributed from a different routing protocol), the router-id is set to the local 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 routing table entry, and the metric is computed as specified in Section 3.5.2.

3.7.1. Periodic Updates

Every Babel speaker periodically advertises all of its selected routes on all of its interfaces, including any recently retracted routes. 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.

3.7.2. Triggered Updates

In addition to the 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, a node **MUST** send a triggered update in a timely manner whenever it changes the selected router-id for a given destination. Additionally, it **SHOULD** make a reasonable attempt at ensuring that all 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 acknowledgement 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 acknowledgements 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, 5 for wireless and 2 for wired links).

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, either due to a received update or because a link cost has changed. 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, 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 garbage-collection timer is not reset when a retraction 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 is not sent on this particular interface when the advertised route was learnt from a neighbour over the same interface.

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).

3.8. Explicit Route Requests

In normal operation, a node's routing table is populated by the regular and triggered updates sent by its neighbours. Under some circumstances, however, a node sends explicit requests 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 a neighbour.

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 update.

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 prefix (prefix, plen), it checks its route table for a selected route to this exact prefix. If such a route exists, it **MUST** send an update; if such a route does not, it **MUST** send a retraction for that prefix.

When a node receives a wildcard route request, it **SHOULD** send a full routing table dump.

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 routing 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 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 and sends an update. A node **MUST NOT** increase its sequence number by more than 1 in response to a seqno request.

Otherwise, if the requested router-id is not its own, the received request's hop count is 2 or more, and the node has a route (not necessarily a feasible one) for the requested prefix that does not use the requestor as a next hop, the node **MUST** forward the request if it has a feasible route to the requested prefix and it is advertising this prefix to neighbours, and **SHOULD** forward the request if it has a (not necessarily feasible) route to the requested prefix. It does so by decreasing the hop count and sending the request in a unicast packet destined to a neighbour that advertises the given prefix and that is not the neighbour from which the request was received.

A node SHOULD maintain a list of recently forwarded requests and forward the reply (an update with a sufficiently large seqno) in a timely manner. A node SHOULD compare every incoming request against its list of recently forwarded requests and avoid forwarding it if it is redundant.

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 for 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 request MUST NOT be forwarded to a multicast address, and it MUST NOT be forwarded to multiple neighbours.

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.

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 attached 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. (Note that, due to duplicate suppression, only a small number of such requests will 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. Under heavy packet loss, however, all such requests might be lost; in that case, the second 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 retract the route.

In order to keep routes from spuriously expiring because they have become unfeasible, a node SHOULD send a unicast seqno request whenever it receives an unfeasible update for a route that is currently selected. The requested sequence number is computed from the source table as 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 should never trigger: 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 retractions in response to non-wildcard route requests (Section 3.8.1.1), this will usually result in either the route being refreshed or a retraction being received.

3.8.2.4. Acquiring New Neighbours

In order to speed up convergence after a mobility event, a node MAY send a unicast wildcard request after acquiring a new neighbour. Additionally, a node MAY send a small number of multicast wildcard requests shortly after booting. Note that doing that carelessly can

cause serious congestion when a whole network is rebooted, especially on link layers with high per-packet overhead (e.g., IEEE 802.11).

4. Protocol Encoding

A Babel packet is 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. Babel packets MUST NOT be sent as IPv6 Jumbograms.

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/IPv4, 48 octets for UDP/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.

In order to avoid global synchronisation of a Babel network and to aggregate multiple TLVs into large packets, a Babel node MUST buffer every TLV and delay sending a UDP packet by a small, randomly chosen delay [JITTER]. In order to allow accurate computation of packet loss rates, this delay MUST NOT be larger than half the advertised Hello interval.

4.1. Data Types

4.1.1. 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.

4.1.2. Router-Id

A router-id is an arbitrary 8-octet. A router-id MUST NOT consist of either all zeroes or all ones. Router-ids SHOULD be assigned in modified EUI-64 format [ADDRARCH].

4.1.3. Address

Since the bulk of the protocol is taken by addresses, multiple ways of encoding addresses are defined. Additionally, a common subnet prefix may be omitted when multiple addresses are sent in a single packet -- this is known as address compression [PACKETBB].

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. The value is 8 octets long, a prefix of fe80::/64 is implied.

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

4.1.4. 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.

```

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
+-----+-----+-----+-----+-----+-----+-----+-----+
|   Magic   |   Version   |   Body length   |
+-----+-----+-----+-----+-----+-----+-----+
| Packet Body ...
+-----+-----+-----+-----+-----+-----+-----+

```

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.

Body The packet body; a sequence of TLVs.

Any data following the body MUST be silently ignored.

4.3. TLV Format

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

```

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      |      Payload...      |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Fields :

Type The type of the TLV.

Length The length of the body, exclusive of the Type and Length fields. If the body is longer than the expected length of a given type of TLV, any extra data MUST be silently ignored.

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 TLV length. If a TLV has a self-terminating format, then it MAY allow a sequence of sub-TLVs to follow the body.

Sub-TLVs have the same structure as TLVs. With the exception of PAD1, all TLVs have the following structure:

```

      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      |      Body...      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

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, called the mandatory 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 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, see Section 4.6.7, Section 4.6.8, and Section 4.6.9).

4.5. Parser state

Babel uses a stateful parser: a TLV may refer to data from a previous TLV. Babel's 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 an Update TLV with flag 80 hexadecimal 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 the 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 both the Router-ID TLV (Section 4.6.7) and the Update TLV with flag 40 hexadecimal set.

Since the parser state is separate from the bulk of Babel's state, and for correct parsing must be identical across implementations, it is updated before checking for mandatory TLVs: parsing a TLV updates

the parser state even if the TLV is otherwise ignored due to an unknown mandatory sub-TLV.

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.

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, exclusive of the Type and Length fields.

MBZ Set to 0 on transmission.

This TLV is silently ignored on reception.

4.6.3. Acknowledgement 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 = 2   |   Length   |   Reserved   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|               Nonce               |   Interval   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```


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

Fields :

Type	Set to 2 to indicate an Acknowledgement Request TLV.
Length	The length of the body, exclusive of the Type and Length fields.
Reserved	Sent as 0 and MUST be ignored on reception.
Nonce	An arbitrary value that will be echoed in the receiver's Acknowledgement 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 acknowledgement before this time has elapsed (with a margin allowing for propagation time).

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

4.6.4. Acknowledgement

```

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 = 3   |   Length   |               Nonce               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

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

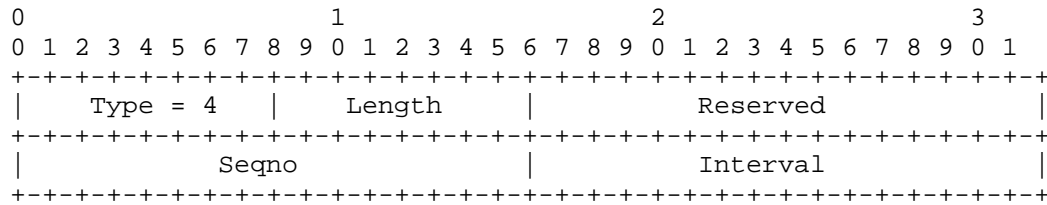
Fields :

Type	Set to 3 to indicate an Acknowledgement TLV.
Length	The length of the body, exclusive of the Type and Length fields.
Nonce	Set to the Nonce value of the Acknowledgement Request that prompted this Acknowledgement.

Since nonce 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



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

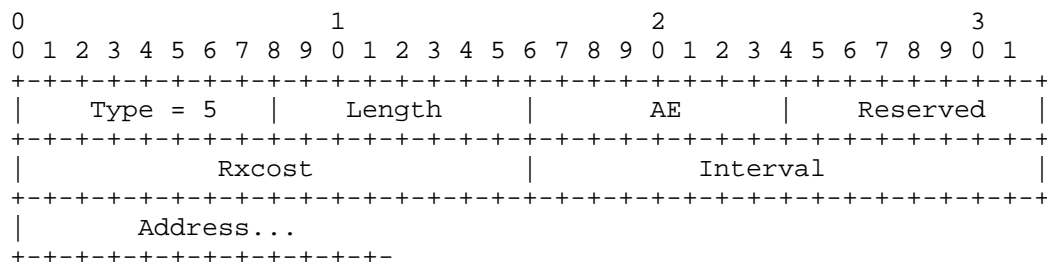
Fields :

- Type Set to 4 to indicate a Hello TLV.
- Length The length of the body, exclusive of the Type and Length fields.
- Reserved Sent as 0 and MUST be ignored on reception.
- Seqno The value of the sending node's Hello seqno for this interface.
- Interval An upper bound, expressed in centiseconds, on the time after which the sending node will send a new Hello TLV. This MUST NOT be 0.

Since there is a single seqno counter for all the Hellos sent by a given node over a given interface, this TLV MUST be sent to a multicast 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



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

Fields :

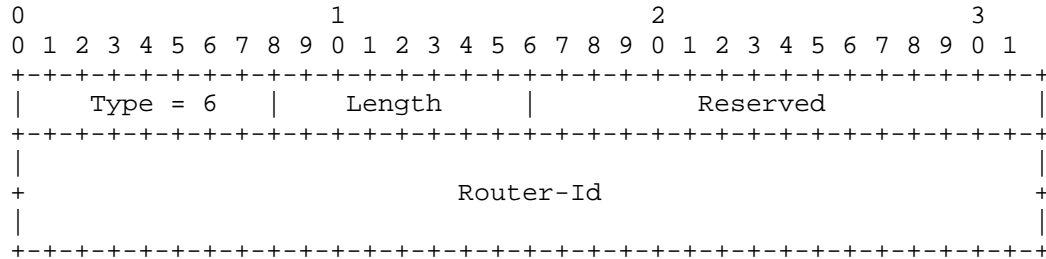
Type	Set to 5 to indicate an IHU TLV.
Length	The length of the body, exclusive of the Type and Length fields.
AE	The encoding of the Address field. This should be 1 or 3 in most cases. As an optimisation, it MAY be 0 if the TLV is sent to a unicast address, if the association is over a point-to-point 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 it SHOULD 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 for 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. 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.

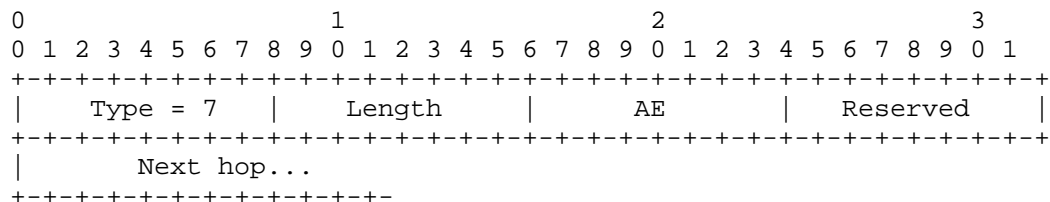
Length The length of the body, exclusive of the Type and Length fields.

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



A Next Hop TLV establishes a next-hop address for a given address family (IPv4 or IPv6) that is implied by subsequent Update TLVs. This TLV sets up the next-hop for subsequent Update TLVs even if it is 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, exclusive of the Type and Length fields.

AE The encoding of the Address field. This SHOULD be 1 or 3 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 that case, the next hop 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

```

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 = 8   |   Length   |   AE   |   Flags   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|   Plen   |   Omitted   |   Interval   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|   Seqno   |   Metric   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|   Prefix...
+-----+-----+-----+-----+-----+-----+-----+-----+

```

An Update TLV advertises or retracts a route. As an optimisation, this can also have the side effect of establishing a new implied router-id and a new default prefix.

Fields :

Type Set to 8 to indicate an Update TLV.

Length The length of the body, 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). Every node MUST be able to interpret the flags with values 80 and 40 hexadecimal; unknown flags MUST be silently ignored.
Plen	The length of the advertised prefix.
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 with the flag with value 80 hexadecimal 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 and SHOULD NOT be less than 10. The receiving node will use this value to compute a hold time for this routing 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:

- o if the bit with value 80 hexadecimal is 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 if the bit with value 40 hexadecimal is 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.

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 flag 80 hexadecimal set and the same address encoding, even if it was ignored due to an unknown mandatory sub-TLV;
- o the next $(\text{Plen}/8 - \text{Omitted})$ rounded upwards octets are taken from the Prefix field;
- 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 bit with value 40 hexadecimal is set in the Flags field, computed as described above. Otherwise, it is taken either from the preceding Router-Id packet, or the preceding Update packet with flag 40 hexadecimal set, whichever comes last, even if that TLV is otherwise ignored due to an unknown mandatory sub-TLV.

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 the metric field is FFFF hexadecimal, this TLV specifies a retraction. In that case, the current router-id and the Seqno are not used. AE MAY then be 0, in which case this Update retracts all of the routes previously advertised on this interface.

Update TLVs with an unknown value for 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 routing table dump.

Fields :

Type	Set to 9 to indicate a Route Request TLV.
Length	The length of the body, 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 routing table dump (a wildcard request).
Plen	The length of the requested prefix.
Prefix	The prefix being requested. This field's size is Plen/8 rounded upwards.

A Request TLV prompts the receiving node to send an update message for the prefix specified by the AE, Plen, and Prefix fields, or a full dump of its routing table if AE is 0 (in which case Plen MUST be 0 and Prefix is of length 0). A Request may be sent to a unicast address if it is destined to a single node, or to a multicast address if the request is destined to all of the neighbours of the sending interface.

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

4.6.11. Seqno 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 = 10   |   Length   |   AE   |   Plen   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|               Seqno        |   Hop Count   |   Reserved   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|
+               Router-Id               +
|
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Prefix...
+---+---+---+---+---+---+---+---+

```

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 message.
Length	The length of the body, exclusive of the Type and Length fields.
AE	The encoding of the Prefix field. This MUST NOT be 0.
Plen	The length 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 an 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
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.

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 Set to 0 on transmission.

This sub-TLV is silently ignored on reception.

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 ff02:0:0:0:0:0:1:6 and the IPv4 multicast group 224.0.0.111 for use by the Babel protocol.

6. Security Considerations

As defined in this document, Babel is a completely insecure protocol. Any attacker can attract data traffic by advertising routes with a low metric. This particular issue can be solved either by lower-layer security mechanisms (e.g., IPsec or link-layer security), or by appending a cryptographic key to Babel packets; the provision of ignoring any data contained within a Babel packet beyond the body length declared by the header is designed for just such a purpose.

The information that a Babel node announces to the whole routing domain is often sufficient to determine a mobile node's physical location with reasonable precision. The privacy issues that this causes can be mitigated somewhat by using randomly chosen router-ids and randomly chosen IP addresses, and changing them periodically.

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 provides protection against spoofed Babel packets being sent from the global Internet. No such natural protection exists when Babel packets are carried over IPv4.

7. References

7.1. Normative References

- [ADDRARCH] Hinden, R. and S. Deering, "IP Version 6 Addressing Architecture", RFC 4291, February 2006.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", RFC 2119, March 1997.

7.2. Informative References

- [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.
- [IS-IS] "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.
- [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 gives a few examples of such strategies.

The sample implementation of Babel maintains statistics about the last 16 received Hello TLVs (Appendix A.1), 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 (Appendix A.3.1).

A.1. Maintaining Hello History

For each neighbour, the sample implementation of Babel maintains a Hello history and an expected sequence number. The Hello history is a vector of 16 bits, where a 1 value represents a received Hello, and

a 0 value a missed 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 from a neighbour, a node compares the received sequence number nr 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 associated neighbour table entry is flushed;
- 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 our 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 neighbour's Hello history, resets the neighbour's Hello timer, and sets ne to $(nr + 1)$. It then resets the neighbour's Hello timer to 1.5 times the value advertised in the received Hello (the extra margin allows for the delay due to jitter).

Whenever the Hello timer associated to a neighbour expires, the local node adds a 0 bit to this neighbour's Hello history, and increments the expected Hello number. If the Hello history is empty (it contains 0 bits only), the neighbour entry is flushed; otherwise, it resets the neighbour's Hello timer to the value advertised in the last Hello received from this neighbour (no extra margin is necessary in this case).

A.2. Cost Computation

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 $K \geq$

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 K ; otherwise, the link is assumed to be down, and the $rxcost$ is set to infinity.

The cost of such a link is defined as

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

A.2.2. ETX

The Estimated Transmission Cost metric [ETX] estimates the number of times that a unicast frame will be retransmitted by the IEEE 802.11 MAC, assuming infinite persistence.

A node uses a neighbour's Hello history to compute an estimate, written β , of the probability that a Hello TLV is successfully received. The $rxcost$ is defined as $256/\beta$.

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

$$cost = 256/(\alpha * \beta)$$

or, equivalently,

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

A.3. Metric Computation

A.3.1. Additive Metrics

The simplest approach for obtaining a monotonic, isotonic metric is to define the metric of a route as the sum of the costs of the component links. More formally, if a neighbour advertises a route with metric m over a link with cost c , then the resulting route has metric $M(c, m) = c + m$.

A multiplicative metric can be converted to an additive one by taking the logarithm (in some suitable base) of the link costs.

A.3.2. External Sources of Willingness

A node may want to vary its willingness to forward packets by taking into account information that is external to the Babel protocol, such as the monetary cost of a link, the node's battery status, CPU load,

etc. This can be done by adding to every route's metric a value k that depends on the external data. For example, if a battery-powered node receives an update with metric m over a link with cost c , it might compute a metric $M(c, m) = k + c + m$, where k depends on the battery status.

In order to preserve strict monotonicity (Section 3.5.2), the value k must be greater than $-c$.

Appendix B. Constants

The choice of time constants is a trade-off between fast detection of mobility events and protocol overhead. Two implementations of Babel with different time constants will interoperate, although the resulting convergence time will most likely be dictated by the slowest of the two implementations.

Experience with the sample implementation of Babel indicates that the Hello interval is the most important time constant: a mobility event is detected within 1.5 to 3 Hello intervals. Due to Babel's reliance on triggered updates and explicit requests, the Update interval only has an effect on the time it takes for accurate metrics to be propagated after variations in link costs too small to trigger an unscheduled update.

At the time of writing, the sample implementation of Babel uses the following default values:

Hello Interval: 4 seconds on wireless links, 20 seconds on wired links.

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

Update Interval: 4 times the Hello interval.

IHU Hold Time: 3.5 times the advertised IHU interval.

Route Expiry Time: 3.5 times the advertised update interval.

Source GC time: 3 minutes.

The amount of jitter applied to a packet depends on whether it contains any urgent TLVs or not. Urgent triggered updates and urgent requests are delayed by no more than 200ms; other TLVs are delayed by no more than one-half the Hello interval.

Appendix C. 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 the mandatory bit set or not);
- o defining new AEs;
- o using the packet trailer.

New versions of the Babel protocol should only be defined 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, or by adding a sub-TLV to the Update TLV.

The two encodings are treated differently by implementations that do not understand the extension. In the case of a new TLV, the whole unknown TLV is ignored by an implementation of the original protocol, while in the case of a new sub-TLV, the TLV is parsed and acted upon, and the unknown sub-TLV is silently ignored. Therefore, a sub-TLV should be used by extensions that extend the Update in a compatible manner (the extension data may be silently ignored), while 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).

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 often 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 similar issue arises with Update TLVs with unknown AEs establishing a new router-id (flag 40 hexadecimal). Therefore, defining new AEs must be done with care if compatibility with unextended implementations is required.

The packet trailer -- the space after the declared length of the packet but within the payload of the UDP datagram -- was originally

intended to carry a cryptographic signature. However, at this time no extension has used it, and therefore we refrain from making any recommendations about its use due to the lack of implementation experience.

Appendix D. Simplified Implementations

Babel is a fairly economic protocol. Route updates take between 12 and 40 octets per destination, depending on how successful compression is; in a double-stack mesh network, an average of less than 24 octets 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 routing table and the associated source table.

Babel is also a reasonably simple protocol. The sample implementation consists of less than 8000 lines of C code, and it compiles to less than 60 kB of text on a 32-bit CISC architecture.

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

A parasitic implementation is one that uses a Babel network for routing its packets but does not announce any of the routes that it has learnt from its neighbours. (This is slightly more than a passive implementation, which doesn't even announce routes to itself.) It may either maintain a full routing table or simply select a gateway amongst any one of its neighbours that announces a default route. Since a parasitic implementation never forwards packets, it cannot possibly participate in a routing loop; hence, it need not evaluate the feasibility condition, and need not maintain a source table.

A parasitic implementation **MUST** answer acknowledgement requests and **MUST** participate in the Hello/IHU protocol. Finally, it **MUST** be able to reply to seqno requests for routes that it announces and **SHOULD** be able to reply to route requests.

Appendix E. Software Availability

The sample implementation of Babel is available from
<<http://www.pps.univ-paris-diderot.fr/~jch/software/babel/>>.

Appendix F. Changes from previous versions

F.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.

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

- o Added security considerations.

F.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 C.
- 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 route 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.

Author's Address

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

Email: jch@irif.fr

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

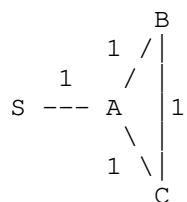
Every node A maintains two pieces of data: its estimated distance to S, written $D(A)$, and its next-hop router to S, written $NH(A)$. Initially, $D(S) = 0$, $D(A)$ is infinite, and $NH(A)$ is undefined.

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.

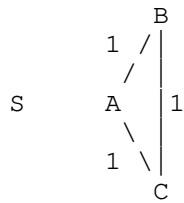
2.3. Transient Loops in Bellman-Ford

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).

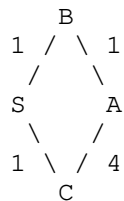
Another simple feasibility condition, used in the Destination-Sequenced Distance-Vector (DSDV) routing protocol [DSDV] and in the Ad hoc On-Demand Distance Vector (AODV) protocol [RFC3561], stems from the following observation: a routing loop can only arise after a router has switched to a route with 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 $D(A) \leq FD(A)$, $D(B) < FD(A)$.

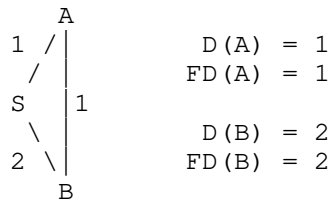
To see that it is strictly less restrictive, consider the following diagram, where A has selected the route through B, and $D(A) = FD(A) = 2$. Since $D(C) = 1 < FD(A)$, the alternate route through C is feasible for A, although its metric $C(A, C) + D(C) = 5$ is larger than that of the currently selected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



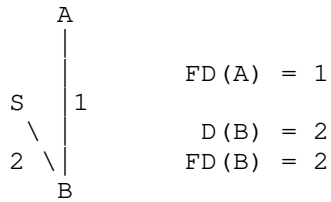
To show that this feasibility condition still guarantees loop-freedom, recall that at the time when A accepts an update from B, the metric $D(B)$ announced by B is no smaller than $FD(B)$; since it is smaller than $FD(A)$, at that point in time $FD(B) < FD(A)$. Since this property is preserved both when 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:

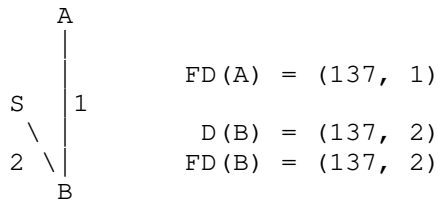


The only route available from A to S, the one that goes through B, is not feasible: A suffers from spurious starvation. At that point, the whole subtree suffering from starvation must be reset, which is essentially what EIGRP does when it performs a global synchronisation of all the routers in the starving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

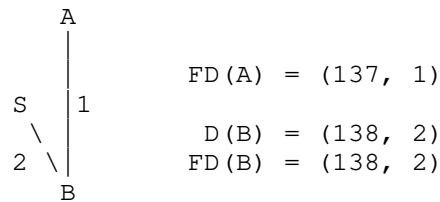
Babel reacts to starvation in a less drastic manner, by using sequenced routes, a technique introduced by DSDV and adopted by AODV. In addition to a metric, every route carries a sequence number, a nondecreasing integer that is propagated unchanged through the network and is only ever incremented by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

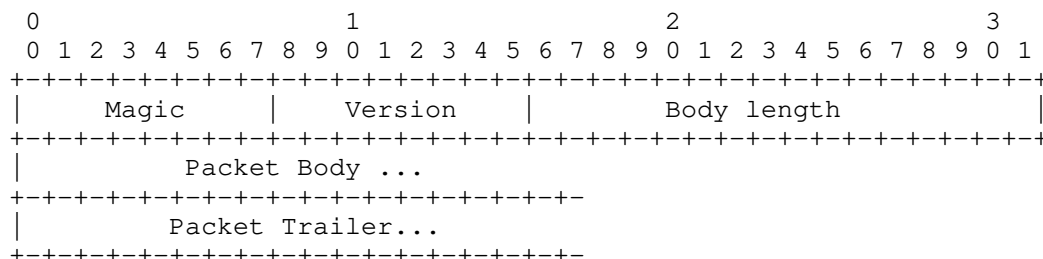
The address family associated to an address encoding is either IPv4 or IPv6; it is undefined for AE 0, IPv4 for AE 1, and IPv6 for AEs 2 and 3.

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).

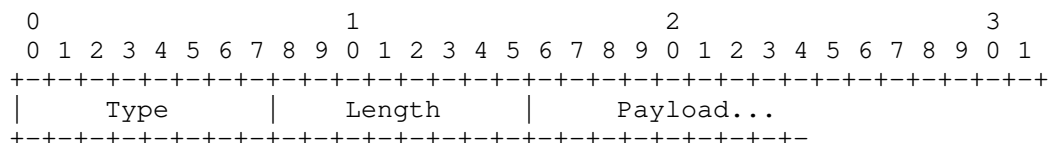
Packet Body The packet body; a sequence of TLVs.

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.

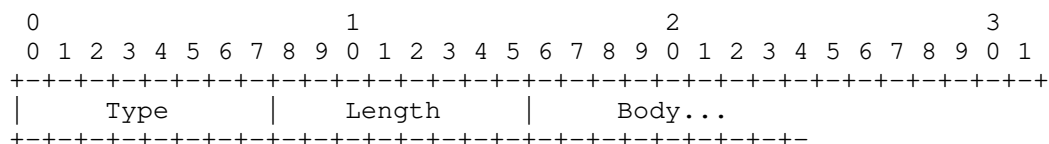
Payload The TLV payload, which consists of a body and, for selected TLV types, an optional list of sub-TLVs.

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

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 = 2										Length										Reserved																			
Opaque										Interval																													

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

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 = 3										Length										Opaque																			

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.

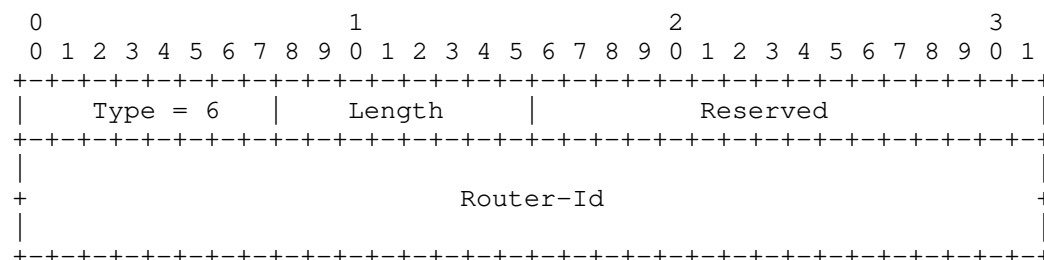
Length	The length of the body in octets, exclusive of the Type and Length fields.
AE	The encoding of the Address field. This should be 1 or 3 in most cases. As an optimisation, it MAY be 0 if the TLV is sent to a unicast address, if the association is over a point-to-point link, or when bidirectional 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.

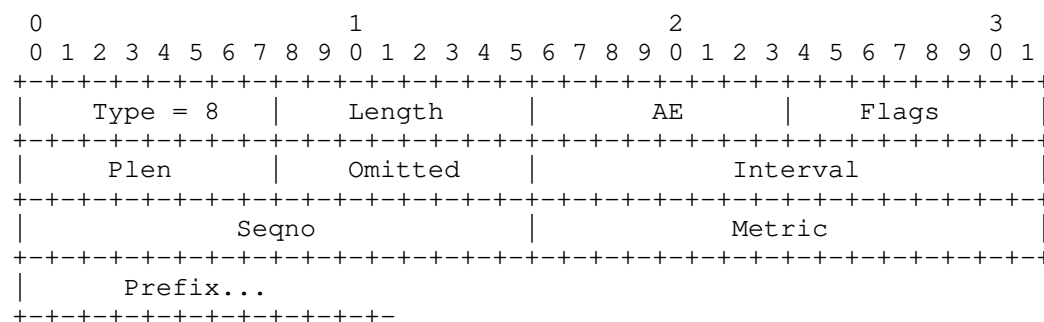
Next hop The next-hop address advertised by subsequent Update TLVs, for this address family.

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.

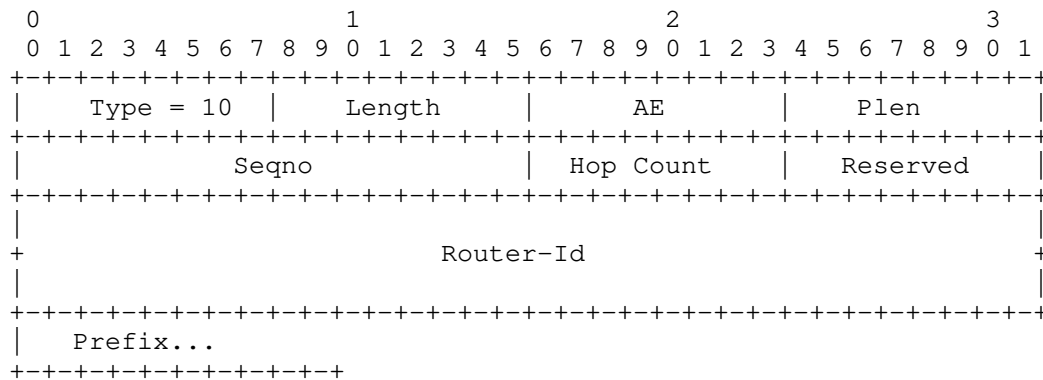
Fields :

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

Babel routing protocol
Internet-Draft
Intended status: Informational
Expires: September 14, 2017

B. Stark
AT&T
March 13, 2017

Babel Information Model
draft-stark-babel-information-model-01

Abstract

This Babel Information Model can be used to create data models under various data modeling regimes (e.g., YANG). It allows a Babel implementation (via a management protocol such as netconf) to report on its current state and may allow some limited configuration of protocol constants.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on September 14, 2017.

Copyright Notice

Copyright (c) 2017 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
1.1. Requirements Language	2
1.2. Notation	3
2. The Information Model	3
2.1. Definition of babel-information	3
2.2. Definition of babel-constants	4
2.3. Definition of babel-interfaces	5
2.4. Definition of babel-neighbors	5
2.5. Definition of babel-csa	6
2.6. Definition of babel-sources	6
2.7. Definition of babel-routes	7
3. References	8
3.1. Normative References	8
3.2. Informative References	8
Author's Address	8

1. Introduction

Babel is a loop-avoiding distance-vector routing protocol defined in RFC 6126 [RFC6126]. Babel Hashed Message Authentication Code (HMAC) Cryptographic Authentication, defined in RFC 7298 [RFC7298], describes a cryptographic authentication mechanism for the Babel routing protocol. This document describes an information model for Babel (including HMAC) that can be used to create management protocol data models (such as a netconf [RFC6241] YANG data model). Other Babel extensions may be included in this document when they become working group drafts.

Due to the simplicity of the Babel protocol and the fact that it is designed to be used in non-professionally administered environments (such as home networks), most of the information model is focused on reporting status of the Babel protocol, and very little of that is considered mandatory to implement (conditional on a management protocol with Babel support being implemented). 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.

1.1. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC2119].

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 values, and n is the maximum. The symbol * for n means no upper bound.

2. The Information Model

2.1. Definition of babel-information

```
object {
    string          babel-version;
    int             babel-self-router-id;
    [int           babel-self-seqno;]
    string          babel-cost-comp-algorithms<1..*>;
    babel-constants-obj babel-constants;
    babel-interfaces-obj babel-interfaces<0..*>;
    babel-sources-obj  babel-sources<0..*>;
    babel-routes-obj   babel-routes<0..*>;
}babel-information-obj;
```

babel-version: the version of this implementation of the Babel protocol

babel-self-router-id: the router-id used by this instance of the Babel protocol to identify itself

babel-self-seqno: the current sequence number included in route updates for routes originated by this node

babel-cost-comp-algorithm: a set of names of supported cost computation algorithms; possible values include "k-out-of-j", "ETX"

babel-constants: a babel-constants object

babel-interfaces: a set of babel-interface objects

babel-sources: a set of babel-source objects

babel-routes: a set of babel-route objects

2.2. Definition of babel-constants

```
object {  
    int          babel-udp-port;  
    [int         babel-multicast-group;]  
}babel-constants-obj;
```

babel-udp-port: UDP port for sending and listening for Babel messages; MAY be configurable

babel-hello-interval-lossy: Hello Interval default for lossy links in milliseconds; MAY be configurable

babel-hello-interval-lossless: Hello Interval default for lossless links in milliseconds; MAY be configurable

babel-ihu-interval: IHU Interval default as multiples of Hello interval

babel-update-interval: Update Interval default as multiples of Hello interval

babel-ihu-hold-time: IHU Hold Time default as multiples of Hello interval

babel-route-expiry-time: IHU Interval default as multiples of Hello interval

babel-garbage-collection-time: Garbage Collection time default as multiples of Update interval

babel-max-trigger-delay: Maximum delay to wait before sending a triggered update in milliseconds

babel-max-normal-delay: Maximum delay to wait before sending a non-triggered message in milliseconds

babel-ack-limit: Threshold for requesting acknowledgements on an interface (do not request acknowledgements if there are more than this many neighbors on the interface); MAY be configurable

babel-resend-trigger-lossy-limit: Resend limit of triggered updates on lossy links (can this be the same, whether or not acknowledgements are requested?)

babel-resend-trigger-lossless-limit: Resend limit of triggered updates on lossless links (can this be the same, whether or not acknowledgements are requested?)

babel-resend-normal-lossy-limit: Resend limit of normal messages on lossy links

babel-resend-normal-lossless-limit: Resend limit of normal messages on lossless links

2.3. Definition of babel-interfaces

```
object {  
    uri                babel-interface-reference;  
    [int               babel-interface-seqno;]  
    [int               babel-interface-hello-interval;]  
    [int               babel-interface-update-interval;]  
    boolean            babel-request-trigger-ack;  
    boolean            babel-lossy-link;  
    [int               babel-external-cost;]  
    babel-neighbors-obj babel-neighbors<1..*>;  
    [babel-csa-obj     babel-csa<1..*>;]  
}babel-interfaces-obj;
```

babel-interface-reference: reference to an interface object as defined by the data model

babel-interface-seqno: the current sequence number in use for this interface

babel-interface-hello-interval: the current hello interval in use for this interface

babel-interface-update-interval: the current update interval in use for this interface

babel-request-trigger-ack: requests acknowledgement of triggered updates (if number of neighbors less than babel-ack-limit); MAY be configurable

babel-lossy-link: indicates (if true) that the link of this interface is considered lossy; MAY be configurable

babel-external-cost: external input to cost of link of this interface (need to determine how to express this);MUST be configurable if implemented

2.4. Definition of babel-neighbors

```
object {
    some address format  babel-neighbor-address;
    string                babel-hello-history;
    int                  babel-txcost;
    int                  babel-hello-seqno;
    int                  babel-neighbor-ihu-interval;
    [int                 babel-rxcost]
}babel-neighbors-obj;
```

babel-neighbor-address: (IPv4 or v6) address the neighbor sends messages from

babel-hello-history: the Hello history (do we want a human readable format?)

babel-txcost: transmission cost value from the last IHU packet received from this neighbor, or FFFF hexadecimal (infinity) if the IHU hold timer for this neighbor has expired

babel-hello-seqno: expected Hello sequence number

babel-neighbor-ihu-interval: current IHU interval for this neighbor

babel-router-id: router-id of the neighbor

babel-rxcost: reception cost calculated for this neighbor

2.5. Definition of babel-csa

```
object {
    string                placeholder;
}babel-csa-obj;
```

placeholder: this section to be filled in, in the future

2.6. Definition of babel-sources

```
object {
    (prefix, plen)        babel-source-prefix;
    int                   babel-source-router-id;
    int                   babel-source-seqno;
    int                   babel-source-metric;
    [int                  babel-source-garbage-collection-time;]
}babel-sources-obj;
```

babel-source-prefix: Prefix (with prefix length)

babel-source-router-id: router-id of the router originating this prefix

babel-source-seqno: last sequence number used by this source

babel-source-metric: this source's feasibility distance

babel-source-garbage-collection-time: garbage-collection timer for this source

2.7. Definition of babel-routes

```
object {  
    (prefix, plen)    babel-route-prefix;  
    int               babel-route-router-id;  
    int               babel-route-neighbor;  
    int               babel-route-metric;  
    int               babel-route-seqno;  
    ip address        babel-route-next-hop;  
    boolean           babel-route-selected;  
}babel-routes-obj;
```

babel-route-prefix: Prefix (with prefix length) for which this route is advertised

babel-route-router-id: router-id of the router originating this prefix

babel-route-neighbor: neighbor that advertised this route (is this a router-id ?)

babel-route-metric: the metric with which this route was advertised by the neighbor, or FFFF hexadecimal (infinity) for a recently retracted route

babel-route-seqno: the sequence number with which this route was advertised

babel-route-next-hop: the next-hop address of this route

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. References

3.1. 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, <<http://www.rfc-editor.org/info/rfc2119>>.

3.2. Informative References

- [RFC6126] Chroboczek, J., "The Babel Routing Protocol", RFC 6126, DOI 10.17487/RFC6126, April 2011, <<http://www.rfc-editor.org/info/rfc6126>>.
- [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, <<http://www.rfc-editor.org/info/rfc6241>>.
- [RFC7298] Ovsienko, D., "Babel Hashed Message Authentication Code (HMAC) Cryptographic Authentication", RFC 7298, DOI 10.17487/RFC7298, July 2014, <<http://www.rfc-editor.org/info/rfc7298>>.

Author's Address

Barbara Stark
AT&T
Atlanta, GA
US

Email: barbara.stark@att.com

BIER WG
Internet-Draft
Intended status: Standards Track
Expires: December 29, 2017

Z. Zhang
ZTE Corporation
A. Przygienda
Juniper Networks
June 27, 2017

BIER in BABEL
draft-zhang-bier-babel-extensions-01

Abstract

BIER introduces a novel multicast architecture. It does not require a signaling protocol to explicitly build multicast distribution trees, nor does it require intermediate nodes to maintain any per-flow state.

Babel defines a distance-vector routing protocol that operates in a robust and efficient fashion both in wired as well as in wireless mesh networks. This document defines a way to carry necessary BIER signaling information in Babel.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on December 29, 2017.

Copyright Notice

Copyright (c) 2017 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents

carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
2. Terminology	2
3. Advertisement of BIER information	2
3.1. BIER BFR-prefix and BIER sub-TLV	3
3.1.1. BIER sub-TLV	3
3.2. BIER MPLS Encapsulation sub-sub-TLV	3
3.3. Optional BIER sub-domain BSL conversion sub-sub-TLV	4
4. Tree types and tunneling	4
5. Security Considerations	5
6. IANA Considerations	5
7. Acknowledgements	5
8. Normative References	5
Authors' Addresses	6

1. Introduction

[I-D.ietf-bier-architecture] introduces a novel multicast architecture. It does not require a signaling protocol to explicitly build multicast distribution trees, nor does it require intermediate nodes to maintain any per-flow state. All procedures necessary to support BIER are abbreviated by the "BIER architecture" moniker in this document.

[RFC6126] and [I-D.ietf-babel-rfc6126bis] define a distance-vector routing protocol under the name of "Babel". Babel operates in a robust and efficient fashion both in ordinary wired as well as in wireless mesh networks.

2. Terminology

The terminology of this documents follows
[I-D.ietf-bier-architecture], [RFC6126], [RFC7557] and
[I-D.ietf-babel-rfc6126bis].

3. Advertisement of BIER information

In case a router is configured with BIER information, and Babel is the routing protocol used, such a router MAY use Babel protocol to announce the BIER information using the BIER sub-TLV specified below.

3.1. BIER BFR-prefix and BIER sub-TLV

BFR-prefix and according information is carried in a Babel Update TLV per [I-D.ietf-babel-rfc6126bis]. A new sub-TLV is defined to convey further BIER information such as BFR-id, sub-domain-id and BSL. Two sub-sub-TLVs are carried as payload of BIER sub-TLV.

The mandatory bit of BIER sub-TLV should be set to 0. If a router cannot recognize a sub-TLV, the router MUST ignore this unknown sub-TLV.

3.1.1. BIER sub-TLV

The BIER sub-TLV format aligns exactly with the definition and restrictions in [I-D.ietf-bier-isis-extensions] and [I-D.ietf-bier-ospf-bier-extensions]. It is a sub-TLV of Babel update TLV. The prefix MUST NOT be summarized and the according sub-TLV MUST be treated as optional and transitive.

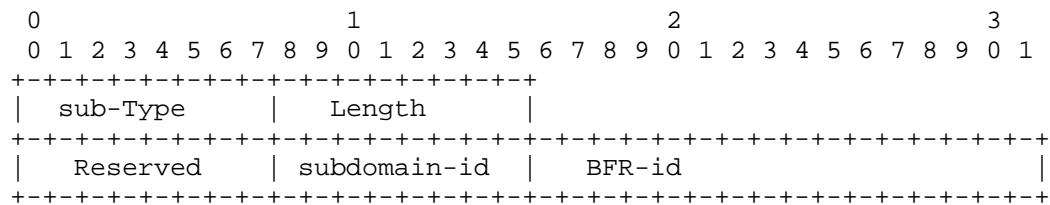


Figure 1: BIER sub-TLV

- o Type: as indicated in IANA section.
- o Length: 1 octet. Include the length of BIER sub-TLV and potential length of the two sub-sub-TLVs.
- o Reserved: MUST be 0 on transmission, ignored on reception. May be used in future versions. 8 bits.
- o subdomain-id: Unique value identifying the BIER sub-domain. 1 octet.
- o BFR-id: A 2 octet field encoding the BFR-id, as documented in [I-D.ietf-bier-architecture]. If no BFR-id has been assigned this field is set to the invalid BFR-id.

3.2. BIER MPLS Encapsulation sub-sub-TLV

The BIER MPLS Encapsulation sub-sub-TLV can be carried by BIER sub-TLV. The format and restrictions are aligned with [I-D.ietf-bier-isis-extensions] and

[I-D.ietf-bier-ospf-bier-extensions]. This sub-sub-TLV carries the information for the BIER MPLS encapsulation including the label range for a specific BSL for a certain <MT,SD> pair.

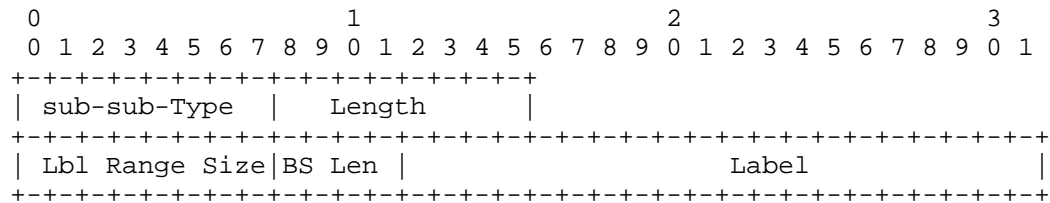


Figure 3: MPLS Encapsulation sub-sub-TLV

- o Type: value of 1 indicating MPLS encapsulation.
- o Length: 1 octet
- o Local BitString Length (BS Len): Encoded bitstring length as per [I-D.ietf-bier-mpls-encapsulation]. 4 bits.
- o Label Range Size: Number of labels in the range for this BIER sub-domain and bitstring length combination, 1 octet.
- o Label: First label of the range, 20 bits. The labels are as defined in [I-D.ietf-bier-mpls-encapsulation].

3.3. Optional BIER sub-domain BSL conversion sub-sub-TLV

This sub-sub-TLV is used to carry the BSL information. Its definition and restrictions are aligned with [I-D.ietf-bier-isis-extensions].

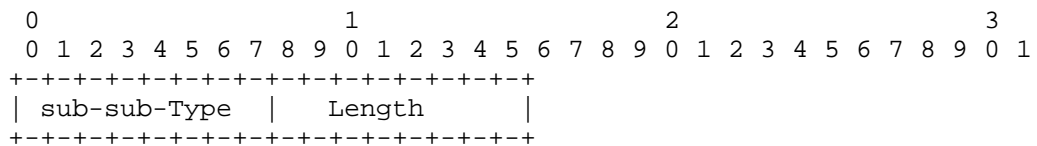


Figure 4: BSL conversion sub-sub-TLV

4. Tree types and tunneling

Since Babel is performing a diffusion computation, support for different tree types is not as natural as with link-state protocols. Hence this specification is assuming that normal Babel reachability computation is performed without further modifications.

BIER architecture does not rely on all routers in a domain performing BFR procedures. How to support tunnels that will allow to tunnel BIER across such routers in Babel is for further study.

5. Security Considerations

TBD

6. IANA Considerations

A new type of Babel update sub-TLV needs to be defined for BIER information advertisement.

7. Acknowledgements

The draft is aligned with the [I-D.ietf-bier-isis-extensions] and [I-D.ietf-bier-ospf-bier-extensions] as far as feasible.

8. Normative References

[I-D.ietf-babel-rfc6126bis]

Chroboczek, J., "The Babel Routing Protocol", draft-ietf-babel-rfc6126bis-02 (work in progress), May 2017.

[I-D.ietf-bier-architecture]

Wijnands, I., Rosen, E., Dolganow, A., Przygienda, T., and S. Aldrin, "Multicast using Bit Index Explicit Replication", draft-ietf-bier-architecture-07 (work in progress), June 2017.

[I-D.ietf-bier-isis-extensions]

Ginsberg, L., Przygienda, T., Aldrin, S., and Z. Zhang, "BIER support via ISIS", draft-ietf-bier-isis-extensions-04 (work in progress), March 2017.

[I-D.ietf-bier-mpls-encapsulation]

Wijnands, I., Rosen, E., Dolganow, A., Tantsura, J., Aldrin, S., and I. Meilik, "Encapsulation for Bit Index Explicit Replication in MPLS and non-MPLS Networks", draft-ietf-bier-mpls-encapsulation-07 (work in progress), June 2017.

[I-D.ietf-bier-ospf-bier-extensions]

Psenak, P., Kumar, N., Wijnands, I., Dolganow, A., Przygienda, T., Zhang, Z., and S. Aldrin, "OSPF Extensions for BIER", draft-ietf-bier-ospf-bier-extensions-06 (work in progress), June 2017.

[RFC6126] Chroboczek, J., "The Babel Routing Protocol", RFC 6126, DOI 10.17487/RFC6126, April 2011, <<http://www.rfc-editor.org/info/rfc6126>>.

[RFC7557] Chroboczek, J., "Extension Mechanism for the Babel Routing Protocol", RFC 7557, DOI 10.17487/RFC7557, May 2015, <<http://www.rfc-editor.org/info/rfc7557>>.

Authors' Addresses

Zheng(Sandy) Zhang
ZTE Corporation
No. 50 Software Ave, Yuhuatai Distinct
Nanjing
China

Email: zhang.zheng@zte.com.cn

Tony Przygienda
Juniper Networks

Email: prz@juniper.net

BIER WG
Internet-Draft
Intended status: Standards Track
Expires: 2 November 2022

Z. Zhang
ZTE Corporation
A. Przygienda
Juniper Networks
1 May 2022

BIER in BABEL
draft-zhang-bier-babel-extensions-07

Abstract

BIER introduces a novel multicast architecture. It does not require a signaling protocol to explicitly build multicast distribution trees, nor does it require intermediate nodes to maintain any per-flow state.

Babel defines a distance-vector routing protocol that operates in a robust and efficient fashion both in wired as well as in wireless mesh networks. This document defines a way to carry necessary BIER signaling information in Babel.

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 2 November 2022.

Copyright Notice

Copyright (c) 2022 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights

and restrictions with respect to this document. Code Components extracted from this document must include Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	2
2. Terminology	2
3. Conventions Used in This Document	3
4. Advertisement of BIER information	3
4.1. BIER BFR-prefix and BIER sub-TLV	3
4.1.1. BIER sub-TLV	3
4.2. BIER MPLS Encapsulation sub-sub-TLV	4
4.3. BIER non-MPLS Encapsulation sub-sub-TLV	5
4.3.1. BIER IPv6 transportation sub-sub-TLV	6
5. Tree types and tunneling	6
6. Security Considerations	6
7. IANA Considerations	6
8. References	6
8.1. Normative References	6
8.2. Informative References	7
Authors' Addresses	8

1. Introduction

[RFC8279] introduces a novel multicast architecture. It does not require a signaling protocol to explicitly build multicast distribution trees, nor does it require intermediate nodes to maintain any per-flow state. All procedures necessary to support BIER are abbreviated by the "BIER architecture" moniker in this document.

[RFC8966] define a distance-vector routing protocol under the name of "Babel". Babel operates in a robust and efficient fashion both in ordinary wired as well as in wireless mesh networks.

2. Terminology

The terminology of this documents follows [RFC8279] and [RFC8966].

3. Conventions Used in This Document

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.

4. Advertisement of BIER information

In case a router is configured with BIER information, and Babel is the routing protocol used, such a router MAY use Babel protocol to announce the BIER information using the BIER sub-TLV specified below.

4.1. BIER BFR-prefix and BIER sub-TLV

BFR-prefix and according information is carried in a Babel Update TLV per [RFC8966]. A new sub-TLV is defined to convey further BIER information such as BFR-id, sub-domain-id and BSL. Two sub-sub-TLVs are carried as payload of BIER sub-TLV.

The mandatory bit of BIER sub-TLV should be set to 0. If a router cannot recognize a sub-TLV, the router MUST ignore this unknown sub-TLV.

4.1.1. BIER sub-TLV

The BIER sub-TLV format aligns exactly with the definition and restrictions in [RFC8401], [RFC8444] and [I-D.ietf-bier-ospfv3-extensions]. It is a sub-TLV of Babel update TLV. The prefix MUST NOT be summarized and the according sub-TLV MUST be treated as optional and transitive.

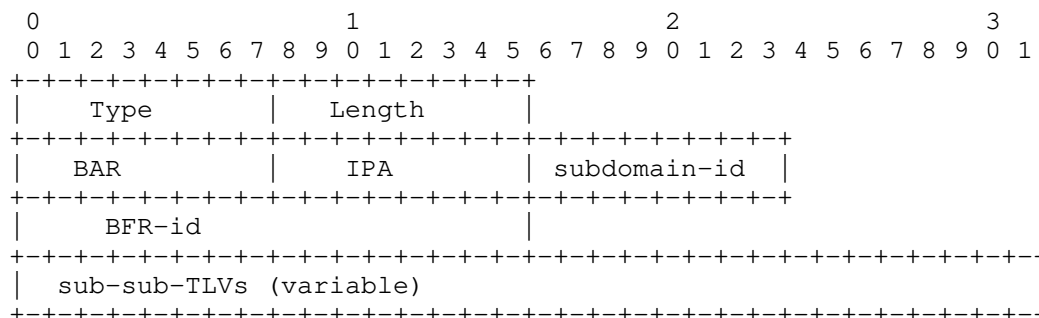


Figure 1: BIER sub-TLV

- * Type: as indicated in IANA section.
- * Length: 1 octet. Include the length of BIER sub-TLV and potential length of the sub-sub-TLVs.
- * BAR: BIER Algorithm. Specifies a BIER-specific algorithm used to calculate underlay paths to reach BFERs. Values are allocated from the "BIER Algorithms" registry. 1 octet.
- * IPA: IGP Algorithm. Specifies an IGP Algorithm to either modify, enhance, or replace the calculation of underlay paths to reach BFERs as defined by the BAR value. Values are from the IGP Algorithm registry. 1 octet.
- * subdomain-id: Unique value identifying the BIER sub-domain. 1 octet.
- * BFR-id: A 2 octet field encoding the BFR-id, as documented in [RFC8279]. If no BFR-id has been assigned this field is set to the invalid BFR-id.

4.2. BIER MPLS Encapsulation sub-sub-TLV

The BIER MPLS Encapsulation sub-sub-TLV can be carried by BIER sub-TLV. The format and restrictions are aligned with [RFC8401], [RFC8444] and [I-D.ietf-bier-ospfv3-extensions]. This sub-sub-TLV carries the information for the BIER MPLS encapsulation including the label range for a specific BSL for a certain <MT,SD> pair.

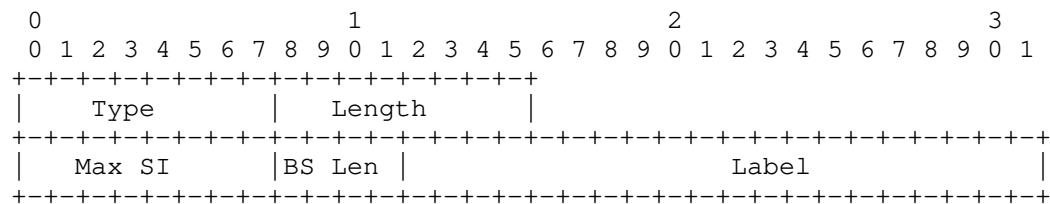


Figure 2: MPLS Encapsulation sub-sub-TLV

- * Type: value of 1 indicating MPLS encapsulation.
- * Length: 1 octet

- * Max SI: Maximum Set Identifier (Section 1 of [RFC8279]) used in the encapsulation for this BIER subdomain for this BitString length, 1 octet. Each SI maps to a single label in the label range. The first label is for SI=0, the second label is for SI=1, etc. If the label associated with the Maximum Set Identifier exceeds the 20-bit range, the sub-sub-TLV MUST be ignored.
- * Local BitString Length (BS Len): Encoded BitString length as per [RFC8296]. 4 bits.
- * Label: First label, 20 bits. The labels are as defined in [RFC8296].

4.3. BIER non-MPLS Encapsulation sub-sub-TLV

The BIER non-MPLS Encapsulation sub-sub-TLV can be carried by BIER sub-TLV. The format and restrictions are aligned with [I-D.ietf-bier-lsr-non-mpls-extensions]. This sub-sub-TLV carries the information for the BIER MPLS encapsulation including the label range for a specific BSL for a certain <MT,SD> pair.

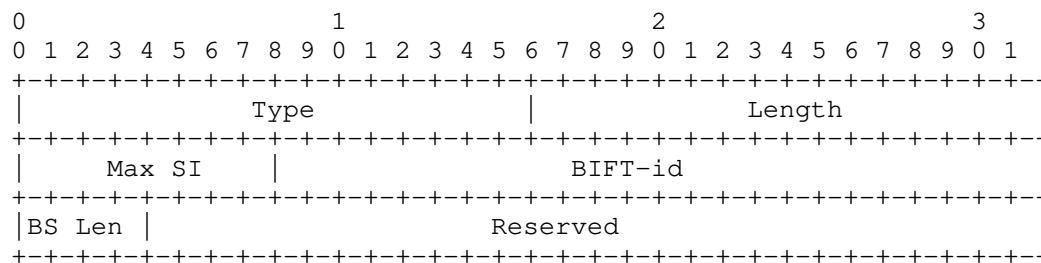


Figure 3: non-MPLS Encapsulation sub-sub-TLV

- * Type: value of 2 indicating non-MPLS encapsulation.
- * Length: 1 octet
- * Max SI: Maximum Set Identifier (Section 1 of [RFC8279]) used in the encapsulation for this BIER subdomain for this BitString length, 1 octet. The first BIFT-id is for SI=0, the second BIFT-id is for SI=1, etc. If the BIFT-id associated with the Maximum Set Identifier exceeds the 20-bit range, the sub-sub-TLV MUST be ignored.
- * BIFT-id: A 3-octet field, where the 20 rightmost bits represent the first BIFT-id in the BIFT-id range. The 4 leftmost bits MUST be ignored. The "BIFT-id range" is the set of 20-bit values

beginning with the BIFT-id and ending with (BIFT-id + (Max SI)). These BIFT-id's are used for BIER forwarding as described in [RFC8279] and [RFC8296].

- * Local BitString Length (BS Len): Encoded BitString length as per [RFC8296]. 4 bits.

4.3.1. BIER IPv6 transportation sub-sub-TLV

The BIER IPv6 transportation sub-sub-TLV can be carried by BIER non-MPLS Encapsulation sub-sub-TLV. The format and restrictions are aligned with [I-D.ietf-bier-bierin6]. A node that requires IPv6 encapsulation MUST advertise the BIER IPv6 transportation sub-sub-TLV according to local configuration or policy in the BIER domain to request other BFRs to always use IPv6 encapsulation.

The format is the same with the definition in section 4.1, [I-D.ietf-bier-bierin6].

5. Tree types and tunneling

Since Babel is performing a diffusion computation, support for different tree types is not as natural as with link-state protocols. Hence this specification is assuming that normal Babel reachability computation is performed without further modifications.

BIER architecture does not rely on all routers in a domain performing BFR procedures. How to support tunnels that will allow to tunnel BIER across such routers in Babel is for further study.

6. Security Considerations

Security considerations discussed in [RFC8296], [RFC8966] apply to this document

7. IANA Considerations

A new type of Babel update sub-TLV needs to be defined for BIER information advertisement.

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>>.

- [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>>.
- [RFC8279] Wijnands, IJ., Ed., Rosen, E., Ed., Dolganow, A., Przygienda, T., and S. Aldrin, "Multicast Using Bit Index Explicit Replication (BIER)", RFC 8279, DOI 10.17487/RFC8279, November 2017, <<https://www.rfc-editor.org/info/rfc8279>>.
- [RFC8296] Wijnands, IJ., Ed., Rosen, E., Ed., Dolganow, A., Tantsura, J., Aldrin, S., and I. Meilik, "Encapsulation for Bit Index Explicit Replication (BIER) in MPLS and Non-MPLS Networks", RFC 8296, DOI 10.17487/RFC8296, January 2018, <<https://www.rfc-editor.org/info/rfc8296>>.
- [RFC8401] Ginsberg, L., Ed., Przygienda, T., Aldrin, S., and Z. Zhang, "Bit Index Explicit Replication (BIER) Support via IS-IS", RFC 8401, DOI 10.17487/RFC8401, June 2018, <<https://www.rfc-editor.org/info/rfc8401>>.
- [RFC8444] Psenak, P., Ed., Kumar, N., Wijnands, IJ., Dolganow, A., Przygienda, T., Zhang, J., and S. Aldrin, "OSPFv2 Extensions for Bit Index Explicit Replication (BIER)", RFC 8444, DOI 10.17487/RFC8444, November 2018, <<https://www.rfc-editor.org/info/rfc8444>>.
- [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>>.

8.2. Informative References

- [I-D.ietf-bier-bierin6]
Zhang, Z., Zhang, Z., Wijnands, I., Mishra, M., Bidgoli, H., and G. Mishra, "Supporting BIER in IPv6 Networks (BIERin6)", Work in Progress, Internet-Draft, draft-ietf-bier-bierin6-04, 2 March 2022, <<https://www.ietf.org/archive/id/draft-ietf-bier-bierin6-04.txt>>.
- [I-D.ietf-bier-lsr-non-mpls-extensions]
Dhanaraj, S., Yan, G., Wijnands, I., Psenak, P., Zhang, Z., and J. Xie, "LSR Extensions for BIER non-MPLS Encapsulation", Work in Progress, Internet-Draft, draft-ietf-bier-lsr-non-mpls-extensions-00, 1 March 2022, <<https://www.ietf.org/archive/id/draft-ietf-bier-lsr-non-mpls-extensions-00.txt>>.

[I-D.ietf-bier-ospfv3-extensions]

Psenak, P., Nainar, N. K., and I. Wijnands, "OSPFv3
Extensions for BIER", Work in Progress, Internet-Draft,
draft-ietf-bier-ospfv3-extensions-05, 19 November 2021,
<[https://www.ietf.org/archive/id/draft-ietf-bier-ospfv3-
extensions-05.txt](https://www.ietf.org/archive/id/draft-ietf-bier-ospfv3-extensions-05.txt)>.

Authors' Addresses

Zheng (Sandy) Zhang
ZTE Corporation
Email: zhang.zheng@zte.com.cn

Tony Przygienda
Juniper Networks
Email: prz@juniper.net