

LoWPAN simple fragment Recovery
draft-thubert-6lowpan-simple-fragment-recovery-02

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with [Section 6 of BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/lid-abstracts.txt>.

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

This Internet-Draft will expire on November 30, 2008.

Copyright Notice

Copyright (C) The IETF Trust (2008).

Abstract

Considering that 6LoWPAN packets can be as large as 2K bytes and that an 802.15.4 frame with security will carry in the order of 80 bytes of effective payload, a packet might end up fragmented into as many as 25 fragments at the 6LoWPAN shim layer. If a single one of those fragments is lost in transmission, all fragments must be resent, further contributing to the congestion that might have caused the initial packet loss. This draft introduces a simple protocol to recover individual fragments between 6LoWPAN endpoints.

Table of Contents

1. Introduction	3
2. Terminology	3
3. Rationale	4
4. Requirements	5
5. Overview	6
6. New Dispatch types and headers	7
6.1. Recoverable Fragment Dispatch type and Header	7
6.2. Fragment Acknowledgement Dispatch type and Header	8
7. Outstanding Fragments Control	8
8. Security Considerations	10
9. IANA Considerations	10
10. Acknowledgments	10
11. References	10
11.1. Normative References	10
11.2. Informative References	10
Author's Address	11
Intellectual Property and Copyright Statements	12

1. Introduction

Considering that 6LoWPAN packets can be as large as 2K bytes and that a 802.15.4 frame with security will carry in the order of 80 bytes of effective payload, a packet might be fragmented into about 25 fragments at the 6LoWPAN shim layer. This level of fragmentation is much higher than that traditionally experienced over the Internet with IPv4 fragments. At the same time, the use of radios increases the probability of transmission loss and Mesh-Under techniques compound that risk over multiple hops.

Past experience with fragmentation has shown that missassociated or lost fragments can lead to poor network behaviour and, eventually, trouble at application layer. The reader might start his research from [[I-D.mathis-frag-harmful](#)] and follow the references. That experience led to the definition of the Path MTU discovery [[RFC1191](#)] protocol that avoids fragmentation over the Internet.

An end-to-end fragment recovery mechanism might be a good complement to a hop-by-hop MAC level recovery with a limited number of retries. This draft introduces a simple protocol to recover individual fragments between 6LoWPAN endpoints. Specifically in the case of UDP, valuable additional information can be found in UDP Usage Guidelines for Application Designers [[I-D.ietf-tsvwg-udp-guidelines](#)].

2. Terminology

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](#)].

Readers are expected to be familiar with all the terms and concepts that are discussed in "IPv6 over Low-Power Wireless Personal Area Networks (6LoWPANs): Overview, Assumptions, Problem Statement, and Goals" [[RFC4919](#)] and "Transmission of IPv6 Packets over IEEE 802.15.4 Networks" [[RFC4944](#)].

ERP

Error Recovery Procedure.

LoWPAN endpoints

The LoWPAN nodes in charge of generating or expanding a 6LoWPAN header from/to a full IPv6 packet. The LoWPAN endpoints are the points where fragmentation and reassembly take place.

3. Rationale

There are a number of usages for large packets in Wireless Sensor Networks. Such usages may not be the most typical or represent the largest amount of traffic over the LoWPAN; however, the associated functionality can be critical enough to justify extra care for ensuring effective transport of large packets across the LoWPAN.

The list of those usages includes:

Towards the LoWPAN node:

Packages of Commands: A number of commands or a full configuration can be packaged as a single message to ensure consistency and enable atomic execution or complete roll back. Until such commands are fully received and interpreted, the intended operation will not take effect.

Firmware update: For example, a new version of the LoWPAN node software is downloaded from a system manager over unicast or multicast services. Such a reflashing operation typically involves updating a large number of similar 6LoWPAN nodes over a relatively short period of time.

From the LoWPAN node:

Waveform captures: A number of consecutive samples are measured at a high rate for a short time and then transferred from a sensor to a gateway or an edge server as a single large report.

Large data packets: Rich data types might require more than one fragment.

Uncontrolled firmware download or waveform upload can easily result in a massive increase of the traffic and saturate the network.

When a fragment is lost in transmission, all fragments are resent, further contributing to the congestion that caused the initial loss, and potentially leading to congestion collapse.

This saturation may lead to excessive radio interference, or random early discard (leaky bucket) in relaying nodes. Additional queueing and memory congestion may result while waiting for a low power next hop to emerge from its sleeping state.

4. Requirements

This paper proposes a method to recover individual fragments between LoWPAN endpoints. The method is designed to fit the following requirements of a LoWPAN (with or without a Mesh-Under routing protocol):

Number of fragments

The recovery mechanism must support highly fragmented packets, with a maximum of 32 fragments per packet.

Minimum acknowledgement overhead

Because the radio is half duplex, and because of silent time spent in the various medium access mechanisms, an acknowledgement consumes roughly as many resources as data fragment.

The recovery mechanism should be able to acknowledge multiple fragments in a single message.

Controlled latency

The recovery mechanism must succeed or give up within the time boundary imposed by the recovery process of the Upper Layer Protocols.

Support for out-of-order fragment delivery

A Mesh-Under load balancing mechanism such as the ISA100 Data Link Layer can introduce out-of-sequence packets. The recovery mechanism must account for packets that appear lost but are actually only delayed over a different path.

Optional flow control

The aggregation of multiple concurrent flows may lead to the saturation of the radio network and congestion collapse.

The recovery mechanism should provide means for controlling the number of fragments in transit over the LoWPAN.

Backward compatibility

A node that implements this draft should be able to communicate with a node that implements [\[RFC4944\]](#). This draft assumes that compatibility information about the remote LoWPAN endpoint is obtained by external means.

5. Overview

Considering that a multi-hop LoWPAN can be a very sensitive environment due to the limited queueing capabilities of a large population of its nodes, this draft recommends a simple and conservative approach to flow control, based on TCP congestion avoidance.

Congestion on the forward path is assumed in case of packet loss, and packet loss is assumed upon time out.

Congestion on the forward path can also be indicated by an Explicit Congestion Notification (ECN) mechanism. Though whether and how ECN [RFC3168] is carried out over the LoWPAN is out of scope, this draft provides a way for the destination endpoint to echo an ECN indication back to the source endpoint in an acknowledgement message as represented in Figure 3 in [Section 6.2](#).

From the standpoint of a source LoWPAN endpoint, an outstanding fragment is a fragment that was sent but for which no explicit acknowledgement was received yet. This means that the fragment might be on the way, received but not yet acknowledged, or the acknowledgement might be on the way back. It is also possible that either the fragment or the acknowledgement was lost on the way.

Because a meshed LoWPAN might deliver frames out of order, it is virtually impossible to differentiate these situations. In other words, from the sender standpoint, all outstanding fragments might still be in the network and contribute to its congestion. There is an assumption, though, that after a certain amount of time, a frame is either received or lost, so it is not causing congestion anymore. This amount of time can be estimated based on the round trip delay between the LoWPAN endpoints. The method detailed in [RFC2988] is recommended for that computation.

The reader is encouraged to read through "Congestion Control Principles" [RFC2914]. Additionally [RFC2309] and [RFC2581] provide deeper information on why this mechanism is needed and how TCP handles Congestion Control. Basically, the goal here is to manage the amount of fragments present in the network; this is achieved by reducing the number of outstanding fragments over a congested path by throttling the sources.

[Section 7](#) describes how the sender decides how many fragments are (re)sent before an acknowledgement is required, and how the sender adapts that number to the network conditions.

6.2. Fragment Acknowledgement Dispatch type and Header

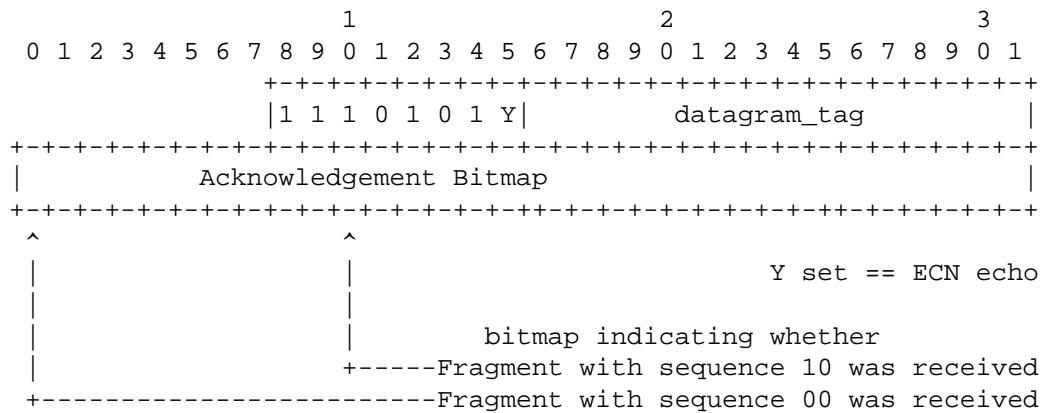


Figure 3: Fragment Acknowledgement Dispatch type and Header

Y bit

When set, the sender indicates that at least one of the acknowledged fragments was received with an Explicit Congestion Notification, indicating that the path followed by the fragments is subject to congestion.

Acknowledgement Bitmap

Each bit in the Bitmap refers to a particular fragment: bit *n* set indicates that fragment with sequence *n* was received, for *n* in [0..31].

All zeroes means that the fragment was dropped because it corresponds to an obsolete `datagram_tag`. This happens if the packet was already reassembled and passed to the network upper layer, or the packet expired and was dropped.

7. Outstanding Fragments Control

A mechanism based on TCP congestion avoidance dictates the maximum number of outstanding fragments.

The maximum number of outstanding fragments for a given packet toward a given LoWPAN endpoint is initially set to a configured value, unless recent history indicates otherwise.

Each time that maximum number of fragments is fully acknowledged,

that number can be incremented by 1. ECN echo and packet loss cause the number to be divided by 2.

The sender transfers a controlled number of fragments and flags the last fragment of a series with an acknowledgement request.

The sender arms a timer to cover the fragment that carries the Acknowledgement request. Upon time out, the sender assumes that all the fragments on the way are received or lost. It divides the maximum number of outstanding fragments by 2 and resets the number of outstanding fragments to 0.

Upon receipt of an Acknowledgement request, the receiver responds with an Acknowledgement containing a bitmap that indicates which fragments were actually received. The bitmap is a 32bit DWORD, which accommodates up to 32 fragments and is sufficient for the 6LoWPAN MTU. For all n in $[0..31]$, bit n is set to 1 in the bitmap to indicate that fragment with sequence n was received, otherwise the bit is set to 0.

The receiver MAY issue unsolicited acknowledgements. An unsolicited acknowledgement enables the sender endpoint to resume sending if it had reached its maximum number of outstanding fragments. Note that acknowledgements might consume precious resources so the use of unsolicited acknowledgements should be configurable and not enabled by default.

The receiver MUST acknowledge a fragment with the acknowledgement request bit set. If any fragment immediately preceding an acknowledgement request is still missing, the receiver MAY intentionally delay its acknowledgement to allow in-transit fragments to arrive. This mechanism might defeat the round trip delay computation so it should be configurable and not enabled by default.

Fragments are sent in a round robin fashion: the sender sends all the fragments for a first time before it retries any lost fragment; lost fragments are retried in sequence, oldest first. This mechanism enables the receiver to acknowledge fragments that were delayed in the network before they are actually retried.

The process must complete within an acceptable time that is within the boundaries of upper layer retries. Additional work is required to define how this is achieved. When the source endpoint decides that a packet should be dropped and the fragmentation process cancelled, it sends a pseudo fragment with the `datagram_offset`, `sequence` and `datagram_size` all set to zero, and no data. Upon reception of this message, the receiver should clean up all resources for the packet associated to the `datagram_tag`.

8. Security Considerations

The process of recovering fragments does not appear to create any opening for new threat.

9. IANA Considerations

Need extensions for formats defined in "Transmission of IPv6 Packets over IEEE 802.15.4 Networks" [[RFC4944](#)].

10. Acknowledgments

The author wishes to thank Jay Werb, Christos Polyzois, Soumitri Kolavennu and Harry Courtice for their contribution and review.

11. References

11.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC2988] Paxson, V. and M. Allman, "Computing TCP's Retransmission Timer", [RFC 2988](#), November 2000.
- [RFC4944] Montenegro, G., Kushalnagar, N., Hui, J., and D. Culler, "Transmission of IPv6 Packets over IEEE 802.15.4 Networks", [RFC 4944](#), September 2007.

11.2. Informative References

- [I-D.ietf-tsvwg-udp-guidelines]
Eggert, L. and G. Fairhurst, "Guidelines for Application Designers on Using Unicast UDP",
[draft-ietf-tsvwg-udp-guidelines-07](#) (work in progress),
May 2008.
- [I-D.mathis-frag-harmful]
Mathis, M., "Fragmentation Considered Very Harmful",
[draft-mathis-frag-harmful-00](#) (work in progress),
July 2004.
- [RFC1191] Mogul, J. and S. Deering, "Path MTU discovery", [RFC 1191](#),
November 1990.

- [RFC2309] Braden, B., Clark, D., Crowcroft, J., Davie, B., Deering, S., Estrin, D., Floyd, S., Jacobson, V., Minshall, G., Partridge, C., Peterson, L., Ramakrishnan, K., Shenker, S., Wroclawski, J., and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet", [RFC 2309](#), April 1998.
- [RFC2581] Allman, M., Paxson, V., and W. Stevens, "TCP Congestion Control", [RFC 2581](#), April 1999.
- [RFC2914] Floyd, S., "Congestion Control Principles", [BCP 41](#), [RFC 2914](#), September 2000.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", [RFC 3168](#), September 2001.
- [RFC4919] Kushalnagar, N., Montenegro, G., and C. Schumacher, "IPv6 over Low-Power Wireless Personal Area Networks (6LoWPANs): Overview, Assumptions, Problem Statement, and Goals", [RFC 4919](#), August 2007.

Author's Address

Pascal Thubert
Cisco Systems
Village d'Entreprises Green Side
400, Avenue de Roumanille
Batiment T3
Biot - Sophia Antipolis 06410
FRANCE

Phone: +33 4 97 23 26 34
Email: pthubert@cisco.com

Full Copyright Statement

Copyright (C) The IETF Trust (2008).

This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights.

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

Intellectual Property

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

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

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

Acknowledgment

Funding for the RFC Editor function is provided by the IETF Administrative Support Activity (IASA).