

Transport Area Working Group
Briscoe
Internet-Draft
Independent

B.

Updates: [6040](#), [2661](#), [2784](#), [3931](#), [4380](#),
2019

July 8,

[7450](#) (if approved)

Intended status: Standards Track

Expires: January 9, 2020

**Propagating Explicit Congestion Notification Across IP Tunnel Headers
Separated by a Shim
draft-ietf-tsvwg-rfc6040update-shim-09**

Abstract

[RFC 6040](#) on "Tunnelling of Explicit Congestion Notification" made the

rules for propagation of ECN consistent for all forms of IP in IP tunnel. This specification updates [RFC 6040](#) to clarify that its scope includes tunnels where two IP headers are separated by at least

one shim header that is not sufficient on its own for wide area packet forwarding. It surveys widely deployed IP tunnelling protocols that use such shim header(s) and updates the specifications

of those that do not mention ECN propagation (L2TPv2, L2TPv3, GRE, Teredo and AMT). This specification also updates [RFC 6040](#) with configuration requirements needed to make any legacy tunnel ingress safe.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

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

Internet-Drafts are draft documents valid for a maximum of six months

and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on January 9, 2020.

Copyright Notice

Copyright (c) 2019 IETF Trust and the persons identified as the document authors. All rights reserved.

Briscoe
1]

Expires January 9, 2020

[Page

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](#) 1. Introduction
- [2](#) 2. Terminology
- [3](#) 3. Scope of [RFC 6040](#)
- [3](#) 3.1. Feasibility of ECN Propagation between Tunnel Headers
- [4](#) 3.2. Desirability of ECN Propagation between Tunnel Headers
- 5 4. Making a non-ECN Tunnel Ingress Safe by Configuration
- [5](#) 5. ECN Propagation and Fragmentation/Reassembly
- [7](#) 6. IP-in-IP Tunnels with Tightly Coupled Shim Headers
- [8](#) 6.1. Specific Updates to Protocols under IETF Change Control
- 11 [6.1.1](#) 6.1.1. L2TP (v2 and v3) ECN Extension
- [11](#) [6.1.2](#) 6.1.2. GRE
- [14](#) [6.1.3](#) 6.1.3. Teredo
- [15](#) [6.1.4](#) 6.1.4. AMT
- [15](#) 7. IANA Considerations
- [17](#) 8. Security Considerations
- [18](#) 9. Comments Solicited
- [18](#) 10. Acknowledgements
- [18](#) 11. References
- [18](#) [11.1](#) 11.1. Normative References
- [18](#) [11.2](#) 11.2. Informative References

[1.](#) Introduction

[RFC 6040](#) on "Tunnelling of Explicit Congestion Notification" [[RFC6040](#)] made the rules for propagation of Explicit Congestion Notification (ECN [[RFC3168](#)]) consistent for all forms of IP in IP tunnel.

A common pattern for many tunnelling protocols is to encapsulate an inner IP header (v4 or v6) with shim header(s) then an outer IP header (v4 or v6). Some of these shim headers are designed as generic encapsulations, so they do not necessarily directly encapsulate an inner IP header. Instead they can encapsulate headers such as link-layer (L2) protocols that in turn often encapsulate IP.

To clear up confusion, this specification clarifies that the scope of [RFC 6040](#) includes any IP-in-IP tunnel, including those with shim header(s) and other encapsulations between the IP headers. Where necessary, it updates the specifications of the relevant encapsulation protocols with the specific text necessary to comply with [RFC 6040](#).

This specification also updates [RFC 6040](#) to state how operators ought to configure a legacy tunnel ingress to avoid unsafe system configurations.

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 [RFC 2119](#) [[RFC2119](#)] when, and only when, they appear in all capitals, as shown here.

This specification uses the terminology defined in [RFC 6040](#) [[RFC6040](#)].

3. Scope of [RFC 6040](#)

In [section 1.1 of RFC 6040](#), its scope is defined as:

"...ECN field processing at encapsulation and decapsulation for any IP-in-IP tunnelling, whether IPsec or non-IPsec tunnels. It applies irrespective of whether IPv4 or IPv6 is used for either the inner or outer headers. ..."

This was intended to include cases where shim header(s) sit between the IP headers. Many tunnelling implementers have interpreted the scope of [RFC 6040](#) as it was intended, but it is ambiguous. Therefore, this specification updates [RFC 6040](#) by adding the following scoping text after the sentences quoted above:

It applies in cases where an outer IP header encapsulates an inner IP header either directly or indirectly by encapsulating other headers that in turn encapsulate (or might encapsulate) an inner IP header.

There is another problem with the scope of [RFC 6040](#). Like many IETF specifications, [RFC 6040](#) is written as a specification that implementations can choose to claim compliance with. This means it does not cover two important cases:

Briscoe
3]

Expires January 9, 2020

[Page

1. those cases where it is infeasible for an implementation to access an inner IP header when adding or removing an outer IP header;
2. those implementations that choose not to propagate ECN between IP headers.

However, the ECN field is a non-optional part of the IP header (v4 and v6). So any implementation that creates an outer IP header has to give the ECN field some value. There is only one safe value a tunnel ingress can use if it does not know whether the egress supports propagation of the ECN field; it has to clear the ECN field in any outer IP header to 0b00.

However, an RFC has no jurisdiction over implementations that choose not to comply with it or cannot comply with it, including all those implementations that pre-dated the RFC. Therefore it would have been

unreasonable to add such a requirement to [RFC 6040](#). Nonetheless, to ensure safe propagation of the ECN field over tunnels, it is reasonable to add requirements on operators, to ensure they configure

their tunnels safely (where possible). Before stating these configuration requirements in [Section 4](#), the factors that determine whether propagating ECN is feasible or desirable will be briefly introduced.

3.1. Feasibility of ECN Propagation between Tunnel Headers

In many cases shim header(s) and an outer IP header are always added to (or removed from) an inner IP packet as part of the same procedure. We call this a tightly coupled shim header. Processing the shim and outer together is often necessary because the shim(s) are not sufficient for packet forwarding in their own right; not unless complemented by an outer header. In these cases it will often

be feasible for an implementation to propagate the ECN field between the IP headers.

In some cases a tunnel adds an outer IP header and a tightly coupled shim header to an inner header that is not an IP header, but that in turn encapsulates an IP header (or might encapsulate an IP header). For instance an inner Ethernet (or other link layer) header might encapsulate an inner IP header as its payload. We call this a tightly coupled shim over an encapsulating header.

Digging to arbitrary depths to find an inner IP header within an encapsulation is strictly a layering violation so it cannot be a required behaviour. Nonetheless, some tunnel endpoints already look within a L2 header for an IP header, for instance to map the Diffserv

codepoint between an encapsulated IP header and an outer IP header

Briscoe
4]

Expires January 9, 2020

[Page

[RFC2983]. In such cases at least, it should be feasible to also (independently) propagate the ECN field between the same IP headers. Thus, access to the ECN field within an encapsulating header can be

a

useful and benign optimization. The guidelines in section 5 of [I-D.ietf-tsvwg-ecn-encap-guidelines] give the conditions for this layering violation to be benign.

3.2. Desirability of ECN Propagation between Tunnel Headers

Developers and network operators are encouraged to implement and deploy tunnel endpoints compliant with RFC 6040 (as updated by the present specification) in order to provide the benefits of wider ECN deployment [RFC8087]. Nonetheless, propagation of ECN between IP headers, whether separated by shim headers or not, has to be

optional

to implement and to use, because:

- o Legacy implementations of tunnels without any ECN support already exist

- o A network might be designed so that there is usually no bottleneck within the tunnel

- o If the tunnel endpoints would have to search within an L2 header to find an encapsulated IP header, it might not be worth the potential performance hit

4. Making a non-ECN Tunnel Ingress Safe by Configuration

Even when no specific attempt has been made to implement propagation of the ECN field at a tunnel ingress, it ought to be possible for

the

operator to render a tunnel ingress safe by configuration. The main safety concern is to disable (clear to zero) the ECN capability in the outer IP header at the ingress if the egress of the tunnel does not implement ECN logic to propagate any ECN markings into the

packet

forwarded beyond the tunnel. Otherwise the non-ECN egress could discard any ECN marking introduced within the tunnel, which would break all the ECN-based control loops that regulate the traffic load over the tunnel.

Therefore this specification updates RFC 6040 by inserting the following text at the end of [section 4.3](#):

"

Whether or not an ingress implementation claims compliance with RFC 6040, RFC 4301 or RFC3168, when the outer tunnel header is IP (v4 or v6), if possible, the operator MUST configure the ingress to zero the outer ECN field in any of the following cases:

Briscoe
5]

Expires January 9, 2020

[Page

- * if it is known that the tunnel egress does not support any of the RFCs that define propagation of the ECN field ([RFC 6040](#), [RFC 4301](#) or the full functionality mode of [RFC 3168](#))
- * or if the behaviour of the egress is not known or an egress with unknown behaviour might be dynamically paired with the ingress.
- * or if an IP header might be encapsulated within a non-IP header that the tunnel ingress is encapsulating, but the ingress does not inspect within the encapsulation.

For the avoidance of doubt, the above only concerns the outer IP header. The ingress MUST NOT alter the ECN field of the arriving IP header that will become the inner IP header.

In order that the network operator can comply with the above safety rules, even if an implementation of a tunnel ingress does not claim to support [RFC 6040](#), [RFC 4301](#) or the full functionality mode of [RFC 3168](#):

- * it MUST NOT treat the former ToS octet (IPv4) or the former Traffic Class octet (IPv6) as a single 8-bit field, as the resulting linkage of ECN and Diffserv field propagation between inner and outer is not consistent with the definition of the 6-bit Diffserv field in [[RFC2474](#)] and [[RFC3260](#)];
- * it SHOULD be able to be configured to zero the ECN field of the outer header.

"

For instance, if a tunnel ingress with no ECN-specific logic had a configuration capability to refer to the last 2 bits of the old ToS Byte of the outer (e.g. with a 0x3 mask) and set them to zero, while also being able to allow the DSCP to be re-mapped independently, that would be sufficient to satisfy both the above implementation requirements.

There might be concern that the above "MUST NOT" makes compliant implementations non-compliant at a stroke. However, by definition it solely applies to equipment that provides Diffserv configuration. Any such Diffserv equipment that is configuring treatment of the former ToS octet (IPv4) or the former Traffic Class octet (IPv6) as a single 8-bit field must have always been non-compliant with the definition of the 6-bit Diffserv field in [[RFC2474](#)] and [[RFC3260](#)].

If a tunnel ingress does not have any ECN logic, copying the ECN field as a side-effect of copying the DSCP is a seriously unsafe bug

Briscoe
6]

Expires January 9, 2020

[Page

that risks breaking the feedback loops that regulate load on a tunnel.

Zeroing the outer ECN field of all packets in all circumstances would

be safe, but it would not be sufficient to claim compliance with [RFC 6040](#) because it would not meet the aim of introducing ECN support to tunnels (see [Section 4.3 of \[RFC6040\]](#)).

5. ECN Propagation and Fragmentation/Reassembly

The following requirements update [RFC6040](#), which omitted handling of the ECN field during fragmentation or reassembly. These changes might alter how many ECN-marked packets are propagated by a tunnel that fragments packets, but this would not raise any backward compatibility issues:

If a tunnel ingress fragments a packet, it MUST set the outer ECN field of all the fragments to the same value as it would have set if it had not fragmented the packet.

As a tunnel egress reassembles sets of outer fragments [\[I-D.ietf-intarea-tunnels\]](#) into packets, it SHOULD propagate CE markings on the basis that a congestion indication on a packet applies to all the octets in the packet. On average, a tunnel egress

SHOULD approximately preserve the number of CE-marked and ECT(1)-marked octets arriving and leaving (counting the size of inner headers, but not encapsulating headers that are being stripped). This process proceeds irrespective of the addresses on the inner headers.

Even if only enough incoming CE-marked octets have arrived for part of the departing packet, the next departing packet SHOULD be immediately CE-marked. This ensures that CE-markings are propagated immediately, rather than held back waiting for more incoming CE-marked octets. Once there are no outstanding CE-marked octets, if only enough incoming ECT(1)-marked octets have arrived for part of the departing packet, the next departing packet SHOULD be immediately marked ECT(1).

For instance, an algorithm for marking departing packets could maintain a pair of counters, the first representing the balance of arriving CE-marked octets minus departing CE-marked octets and the second representing a similar balance of ECT(1)-marked octets. The algorithm:

- o adds the size of every CE-marked or ECT(1)-marked packet that arrives to the appropriate counter;

Briscoe
7]

Expires January 9, 2020

[Page

- o if the CE counter is positive, it CE-marks the next packet to depart and subtracts its size from the CE counter;
- o if the CE counter is negative but the ECT(1) counter is positive, it marks the next packet to depart as ECT(1) and subtracts its size from the ECT((1) counter;
- o (the previous two steps will often leave a negative remainder in the counters, which is deliberate);
- o if neither counter is positive, it marks the next packet to depart as ECT(0);
- o until all the fragments of a packet have arrived, it does not commit any updates to the counters so that, if reassembly fails and the partly reassembled packet has to be discarded, none of the discarded fragments will have updated any of the counters.

During reassembly of outer fragments [[I-D.ietf-intarea-tunnels](#)], if the ECN fields of the outer headers being reassembled into a single packet consist of a mixture of Not-ECT and other ECN codepoints, the packet MUST be discarded.

A tunnel end-point that claims to support the present specification MUST NOT use an approach that results in a significantly different ECN-marking outcome to that defined by the "SHOULD" statements throughout this section. "SHOULD" is only used to allow similar perhaps more efficient approaches that result in approximately the same outcome.

6. IP-in-IP Tunnels with Tightly Coupled Shim Headers

There follows a list of specifications of encapsulations with tightly coupled shim header(s), in rough chronological order. The list is confined to standards track or widely deployed protocols. The list is not necessarily exhaustive so, for the avoidance of doubt, the scope of [RFC 6040](#) is defined in [Section 3](#) and is not limited to this list.

- o PPTP (Point-to-Point Tunneling Protocol) [[RFC2637](#)];
- o L2TP (Layer 2 Tunneling Protocol), specifically L2TPv2 [[RFC2661](#)] and L2TPv3 [[RFC3931](#)], which not only includes all the L2-specific specializations of L2TP, but also derivatives such as the Keyed IPv6 Tunnel [[RFC8159](#)];
- o GRE (Generic Routing Encapsulation) [[RFC2784](#)] and NVGRE (Network Virtualization using GRE) [[RFC7637](#)];

Briscoe
8]

Expires January 9, 2020

[Page

- o GTP (GPRS Tunnelling Protocol), specifically GTPv1 [[GTPv1](#)], GTP v1 User Plane [[GTPv1-U](#)], GTP v2 Control Plane [[GTPv2-C](#)];
- o Teredo [[RFC4380](#)];
- o CAPWAP (Control And Provisioning of Wireless Access Points) [[RFC5415](#)];
- o LISP (Locator/Identifier Separation Protocol) [[RFC6830](#)];
- o AMT (Automatic Multicast Tunneling) [[RFC7450](#)];
- o VXLAN (Virtual eXtensible Local Area Network) [[RFC7348](#)] and VXLAN-GPE [[I-D.ietf-nvo3-vxlan-gpe](#)];
- o The Network Service Header (NSH [[RFC8300](#)]) for Service Function Chaining (SFC);
- o Geneve [[I-D.ietf-nvo3-geneve](#)];
- o GUE (Generic UDP Encapsulation) [[I-D.ietf-intarea-gue](#)];
- o Direct tunnelling of an IP packet within a UDP/IP datagram (see [Section 3.1.11 of RFC8085](#));
- o TCP Encapsulation of IKE and IPsec Packets (see [Section 12.5 of RFC8229](#)).

Some of the listed protocols enable encapsulation of a variety of network layer protocols as inner and/or outer. This specification applies in the cases where there is an inner and outer IP header as described in [Section 3](#). Otherwise [[I-D.ietf-tsvwg-ecn-encap-guidelines](#)] gives guidance on how to design propagation of ECN into other protocols that might encapsulate IP.

Where protocols in the above list need to be updated to specify ECN propagation and they are under IETF change control, update text is given in the following subsections. For those not under IETF control, it is RECOMMENDED that implementations of encapsulation and decapsulation comply with [RFC 6040](#). It is also RECOMMENDED that their specifications are updated to add a requirement to comply with [RFC 6040](#) (as updated by the present document).

PPTP is not under the change control of the IETF, but it has been documented in an informational RFC [[RFC2637](#)]. However, there is no need for the present specification to update PPTP because L2TP has been developed as a standardized replacement.

Briscoe
9]

Expires January 9, 2020

[Page

NVGRE is not under the change control of the IETF, but it has been documented in an informational RFC [[RFC7637](#)]. NVGRE is a specific use-case of GRE (it re-purposes the key field from the initial specification of GRE [[RFC1701](#)] as a Virtual Subnet ID). Therefore the text that updates GRE in [Section 6.1.2](#) below is also intended to update NVGRE.

Although the definition of the various GTP shim headers is under the control of the 3GPP, it is hard to determine whether the 3GPP or the IETF controls standardization of the `_process_` of adding both a GTP and an IP header to an inner IP header. Nonetheless, the present specification is provided so that the 3GPP can refer to it from any of its own specifications of GTP and IP header processing.

The specification of CAPWAP already specifies [RFC 3168](#) ECN propagation and ECN capability negotiation. Without modification the

CAPWAP specification already interworks with the backward compatible updates to [RFC 3168](#) in [RFC 6040](#).

LISP made the ECN propagation procedures in [RFC 3168](#) mandatory from the start. [RFC 3168](#) has since been updated by [RFC 6040](#), but the changes are backwards compatible so there is still no need for LISP tunnel endpoints to negotiate their ECN capabilities.

VXLAN is not under the change control of the IETF but it has been documented in an informational RFC. In contrast, VXLAN-GPE (Generic Protocol Extension) is being documented under IETF change control. It is RECOMMENDED that VXLAN and VXLAN-GPE implementations comply with [RFC 6040](#) when the VXLAN header is inserted between (or removed from between) IP headers. The authors of any future update to these specifications are encouraged to add a requirement to comply with [RFC 6040](#) as updated by the present specification.

The Network Service Header (NSH [[RFC8300](#)]) has been defined as a shim-based encapsulation to identify the Service Function Path (SFP) in the Service Function Chaining (SFC) architecture [[RFC7665](#)]. A proposal has been made for the processing of ECN when handling transport encapsulation [[I-D.ietf-sfc-nsh-ecn-support](#)].

The specifications of Geneve and GUE already refer to [RFC 6040](#) for ECN encapsulation.

[Section 3.1.11 of RFC 8085](#) already explains that a tunnel that encapsulates an IP header within a UDP/IP datagram needs to follow [RFC 6040](#) when propagating the ECN field between inner and outer IP headers. The requirements in [Section 4](#) update [RFC 6040](#), and hence implicitly update the UDP usage guidelines in [RFC 8085](#) to add the important but previously unstated requirement that, if the UDP tunnel

Briscoe
10]

Expires January 9, 2020

[Page

egress does not, or might not, support ECN propagation, a UDP tunnel ingress has to clear the outer IP ECN field to 0b00, e.g. by configuration.

[Section 12.5](#) of TCP Encapsulation of IKE and IPsec Packets [[RFC8229](#)] already recommends the compatibility mode of [RFC 6040](#) in this case, because there is not a one-to-one mapping between inner and outer packets.

[6.1.](#) Specific Updates to Protocols under IETF Change Control

[6.1.1.](#) L2TP (v2 and v3) ECN Extension

The L2TP terminology used here is defined in [[RFC2661](#)] and [[RFC3931](#)].

L2TPv3 [[RFC3931](#)] is used as a shim header between any packet-switched network (PSN) header (e.g. IPv4, IPv6, MPLS) and many types of layer 2 (L2) header. The L2TPv3 shim header encapsulates an L2-specific sub-layer then an L2 header that is likely to contain an inner IP header (v4 or v6). Then this whole stack of headers can be encapsulated optionally within an outer UDP header then an outer PSN header that is typically IP (v4 or v6).

L2TPv2 is used as a shim header between any PSN header and a PPP header, which is in turn likely to encapsulate an IP header.

Even though these shims are rather fat (particularly in the case of L2TPv3), they still fit the definition of a tightly coupled shim header over an encapsulating header ([Section 3.1](#)), because all the headers encapsulating the L2 header are added (or removed) together. L2TPv2 and L2TPv3 are therefore within the scope of [RFC 6040](#), as updated by [Section 3](#) above.

L2TP maintainers are RECOMMENDED to implement the ECN extension to L2TPv2 and L2TPv3 defined in [Section 6.1.1.2](#) below, in order to provide the benefits of ECN [[RFC8087](#)], whenever a node within an L2TP tunnel becomes the bottleneck for an end-to-end traffic flow.

[6.1.1.1.](#) Safe Configuration of a 'Non-ECN' Ingress LCCE

The following text is appended to both [Section 5.3 of \[RFC2661\]](#) and [Section 4.5 of \[RFC3931\]](#) as an update to the base L2TPv2 and L2TPv3 specifications:

The operator of an LCCE that does not support the ECN Extension in [Section 6.1.1.2](#) of RFCXXXX MUST follow the configuration requirements in [Section 4](#) of RFCXXXX to ensure it clears the

outer

IP ECN field to 0b00 when the outer PSN header is IP (v4 or v6).

Briscoe
11]

Expires January 9, 2020

[Page

{RFCXXXX refers to the present document so it will need to be inserted by the RFC Editor}

In particular, for an LCCE implementation that does not support the ECN Extension, this means that configuration of how it propagates the ECN field between inner and outer IP headers MUST be independent of any configuration of the Diffserv extension of L2TP [[RFC3308](#)].

6.1.1.2. ECN Extension for L2TP (v2 or v3)

When the outer PSN header and the payload inside the L2 header are both IP (v4 or v6), to comply with [RFC 6040](#), an LCCE will follow the rules for propagation of the ECN field at ingress and egress in [Section 4 of RFC 6040](#) [[RFC6040](#)].

Before encapsulating any data packets, [RFC 6040](#) requires an ingress LCCE to check that the egress LCCE supports ECN propagation as defined in [RFC 6040](#) or one of its compatible predecessors ([[RFC4301](#)] or the full functionality mode of [[RFC3168](#)]). If the egress supports ECN propagation, the ingress LCCE can use the normal mode of encapsulation (copying the ECN field from inner to outer). Otherwise, the ingress LCCE has to use compatibility mode [[RFC6040](#)] (clearing the outer IP ECN field to 0b00).

An LCCE can determine the remote LCCE's support for ECN either statically (by configuration) or by dynamic discovery during setup of each control connection between the LCCEs, using the Capability AVP defined in [Section 6.1.1.2.1](#) below.

Where the outer PSN header is some protocol other than IP that supports ECN, the appropriate ECN propagation specification will need to be followed, e.g. "Explicit Congestion Marking in MPLS" [[RFC5129](#)]. Where no specification exists for ECN propagation by a particular PSN, [[I-D.ietf-tsvwg-ecn-encap-guidelines](#)] gives general guidance on how to design ECN propagation into a protocol that encapsulates IP.

6.1.1.2.1. LCCE Capability AVP for ECN Capability Negotiation

The LCCE Capability Attribute-Value Pair (AVP) defined here has Attribute Type ZZ. The Attribute Value field for this AVP is a bit-mask with the following 16-bit format:

Briscoe
12]

Expires January 9, 2020

[Page

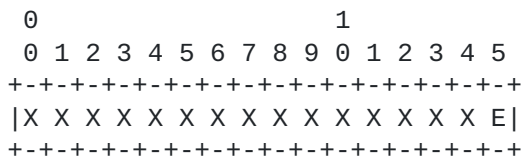


Figure 1: Value Field for the LCCE Capability Attribute

This AVP MAY be present in the following message types: SCCRQ and SCCRP (Start-Control-Connection-Request and Start-Control-Connection-

Reply). This AVP MAY be hidden (the H-bit set to 0 or 1) and is optional (M-bit not set). The length (before hiding) of this AVP MUST be 8 octets. The Vendor ID is the IETF Vendor ID of 0.

Bit 15 of the Value field of the LCCE Capability AVP is defined as the ECN Capability flag (E). When the ECN Capability flag is set to 1, it indicates that the sender supports ECN propagation. When the ECN Capability flag is cleared to zero, or when no LCCE Capability AVP

is present, it indicates that the sender does not support ECN propagation. All the other bits are reserved. They MUST be cleared to zero when sent and ignored when received or forwarded.

An LCCE initiating a control connection will send a Start-Control-Connection-Request (SCCRQ) containing an LCCE Capability AVP with the

ECN Capability flag set to 1. If the tunnel terminator supports ECN, it will return a Start-Control-Connection-Reply (SCCRP) that also includes an LCCE Capability AVP with the ECN Capability flag set to 1. Then, for any sessions created by that control connection, both ends of the tunnel can use the normal mode of [RFC 6040](#), i.e. it can copy the IP ECN field from inner to outer when encapsulating data packets.

If, on the other hand, the tunnel terminator does not support ECN it will ignore the ECN flag in the LCCE Capability AVP and send an SCCRP

to the tunnel initiator without a Capability AVP (or with a Capability AVP but with the ECN Capability flag cleared to zero). The tunnel initiator interprets the absence of the ECN Capability flag in the SCCRP as an indication that the tunnel terminator is incapable of supporting ECN. When encapsulating data packets for any

sessions created by that control connection, the tunnel initiator will then use the compatibility mode of [RFC 6040](#) to clear the ECN field of the outer IP header to 0b00.

If the tunnel terminator does not support this ECN extension, the network operator is still expected to configure it to comply with the

safety provisions set out in [Section 6.1.1.1](#) above, when it acts as an ingress LCCE.

Briscoe
13]

Expires January 9, 2020

[Page

6.1.2. GRE

The GRE terminology used here is defined in [\[RFC2784\]](#). GRE is often used as a tightly coupled shim header between IP headers. Sometimes the GRE shim header encapsulates an L2 header, which might in turn encapsulate an IP header. Therefore GRE is within the scope of [RFC 6040](#) as updated by [Section 3](#) above.

GRE tunnel endpoint maintainers are RECOMMENDED to support [\[RFC6040\]](#) as updated by the present specification, in order to provide the benefits of ECN [\[RFC8087\]](#) whenever a node within a GRE tunnel becomes the bottleneck for an end-to-end IP traffic flow tunnelled over GRE using IP as the delivery protocol (outer header).

GRE itself does not support dynamic set-up and configuration of tunnels. However, control plane protocols such as Mobile IPv4 (MIP4) [\[RFC5944\]](#), Mobile IPv6 (MIP6) [\[RFC6275\]](#), Proxy Mobile IP (PMIP) [\[RFC5845\]](#) and IKEv2 [\[RFC7296\]](#) are sometimes used to set up GRE tunnels dynamically.

When these control protocols set up IP-in-IP or IPSec tunnels, it is likely that they propagate the ECN field as defined in [RFC 6040](#) or one of its compatible predecessors ([RFC 4301](#) or the full functionality mode of [RFC 3168](#)). However, if they use a GRE encapsulation, this presumption is less sound.

Therefore, If the outer delivery protocol is IP (v4 or v6) the operator is obliged to follow the safe configuration requirements in [Section 4](#) above. [Section 6.1.2.1](#) below updates the base GRE specification with this requirement, to emphasize its importance.

Where the delivery protocol is some protocol other than IP that supports ECN, the appropriate ECN propagation specification will need to be followed, e.g Explicit Congestion Marking in MPLS [\[RFC5129\]](#). Where no specification exists for ECN propagation by a particular PSN, [\[I-D.ietf-tsvwg-ecn-encap-guidelines\]](#) gives more general guidance on how to propagate ECN to and from protocols that encapsulate IP.

6.1.2.1. Safe Configuration of a 'Non-ECN' GRE Ingress

The following text is appended to [Section 3 of \[RFC2784\]](#) as an update to the base GRE specification:

The operator of a GRE tunnel ingress MUST follow the configuration requirements in [Section 4](#) of RFCXXXX when the outer delivery protocol is IP (v4 or v6). {RFCXXXX refers to the present

document

so it will need to be inserted by the RFC Editor}

Briscoe
14]

Expires January 9, 2020

[Page

6.1.3. Teredo

Teredo [[RFC4380](#)] provides a way to tunnel IPv6 over an IPv4 network, with a UDP-based shim header between the two.

For Teredo tunnel endpoints to provide the benefits of ECN, the Teredo specification would have to be updated to include negotiation of the ECN capability between Teredo tunnel endpoints. Otherwise it would be unsafe for a Teredo tunnel ingress to copy the ECN field to the IPv6 outer.

It is believed that current implementations do not support propagation of ECN, but that they do safely zero the ECN field in the outer IPv6 header. However the specification does not mention anything about this.

To make existing Teredo deployments safe, it would be possible to add ECN capability negotiation to those that are subject to remote OS update. However, for those implementations not subject to remote OS update, it will not be feasible to require them to be configured correctly, because Teredo tunnel endpoints are generally deployed on hosts.

Therefore, until ECN support is added to the specification of Teredo, the only feasible further safety precaution available here is to update the specification of Teredo implementations with the following text, as a new [section 5.1.3](#):

"5.1.3 Safe 'Non-ECN' Teredo Encapsulation

A Teredo tunnel ingress implementation that does not support ECN propagation as defined in [RFC 6040](#) or one of its compatible predecessors ([RFC 4301](#) or the full functionality mode of [RFC 3168](#)) MUST zero the ECN field in the outer IPv6 header."

6.1.4. AMT

Automatic Multicast Tunneling (AMT [[RFC7450](#)]) is a tightly coupled shim header that encapsulates an IP packet and is itself encapsulated within a UDP/IP datagram. Therefore AMT is within the scope of [RFC 6040](#) as updated by [Section 3](#) above.

AMT tunnel endpoint maintainers are RECOMMENDED to support [[RFC6040](#)] as updated by the present specification, in order to provide the benefits of ECN [[RFC8087](#)] whenever a node within an AMT tunnel becomes the bottleneck for an IP traffic flow tunnelled over AMT.

Briscoe
15]

Expires January 9, 2020

[Page

To comply with [RFC 6040](#), an AMT relay and gateway will follow the rules for propagation of the ECN field at ingress and egress respectively, as described in [Section 4 of RFC 6040](#) [[RFC6040](#)].

Before encapsulating any data packets, [RFC 6040](#) requires an ingress AMT relay to check that the egress AMT gateway supports ECN propagation as defined in [RFC 6040](#) or one of its compatible predecessors ([RFC 4301](#) or the full functionality mode of [RFC 3168](#)). If the egress gateway supports ECN, the ingress relay can use the normal mode of encapsulation (copying the IP ECN field from inner to outer). Otherwise, the ingress relay has to use compatibility mode, which means it has to clear the outer ECN field to zero [[RFC6040](#)].

An AMT tunnel is created dynamically (not manually), so the relay will need to determine the remote gateway's support for ECN using the ECN capability declaration defined in [Section 6.1.4.2](#) below.

6.1.4.1. Safe Configuration of a 'Non-ECN' Ingress AMT Relay

The following text is appended to [Section 4.2.2 of \[RFC7450\]](#) as an update to the AMT specification:

The operator of an AMT relay that does not support [RFC 6040](#) or one of its compatible predecessors ([RFC 4301](#) or the full functionality mode of [RFC 3168](#)) MUST follow the configuration requirements in [Section 4](#) of RFCXXXX to ensure it clears the outer IP ECN field to zero. {RFCXXXX refers to the present document so it will need to be inserted by the RFC Editor}

6.1.4.2. ECN Capability Declaration of an AMT Gateway

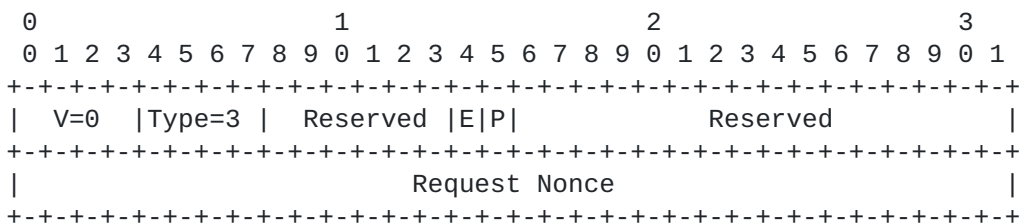


Figure 2: Updated AMT Request Message Format

Bit 14 of the AMT Request Message counting from 0 (or bit 7 of the Reserved field counting from 1) is defined here as the AMT Gateway ECN Capability flag (E), as shown in Figure 2. The definitions of all other fields in the AMT Request Message are unchanged from [RFC 7450](#).

Briscoe
16]

Expires January 9, 2020

[Page

When the E flag is set to 1, it indicates that the sender of the message supports [RFC 6040](#) ECN propagation. When it is cleared to zero, it indicates the sender of the message does not support [RFC 6040](#) ECN propagation. An AMT gateway "that supports [RFC 6040](#) ECN propagation" means one that propagates the ECN field to the forwarded data packet based on the combination of arriving inner and outer ECN fields, as defined in [Section 4 of RFC 6040](#).

The other bits of the Reserved field remain reserved. They will continue to be cleared to zero when sent and ignored when either received or forwarded, as specified in [Section 5.1.3.3. of RFC 7450](#).

An AMT gateway that does not support [RFC 6040](#) MUST NOT set the E flag of its Request Message to 1.

An AMT gateway that supports [RFC 6040](#) ECN propagation MUST set the E flag of its Relay Discovery Message to 1.

The action of the corresponding AMT relay that receives a Request message with the E flag set to 1 depends on whether the relay itself supports [RFC 6040](#) ECN propagation:

- o If the relay supports [RFC 6040](#) ECN propagation, it will store the ECN capability of the gateway along with its address. Then whenever it tunnels datagrams towards this gateway, it MUST use the normal mode of [RFC 6040](#) to propagate the ECN field when encapsulating datagrams (i.e. it copies the IP ECN field from inner to outer).
- o If the discovered AMT relay does not support [RFC 6040](#) ECN propagation, it will ignore the E flag in the Reserved field, as per [section 5.1.3.3. of RFC 7450](#).

If the AMT relay does not support [RFC 6040](#) ECN propagation, the network operator is still expected to configure it to comply with the safety provisions set out in [Section 6.1.4.1](#) above.

7. IANA Considerations

IANA is requested to assign the following L2TP Control Message Attribute Value Pair:

| Attribute Type | Description | Reference |
|----------------|----------------|-----------|
| ZZ | ECN Capability | RFCXXXX |

Briscoe
17]

Expires January 9, 2020

[Page

[TO BE REMOVED: This registration should take place at the following location: <https://www.iana.org/assignments/l2tp-parameters/l2tp-parameters.xhtml>]

8. Security Considerations

The Security Considerations in [[RFC6040](#)] and [[I-D.ietf-tsvwg-ecn-encap-guidelines](#)] apply equally to the scope defined for the present specification.

9. Comments Solicited

Comments and questions are encouraged and very welcome. They can be addressed to the IETF Transport Area working group mailing list <tsvwg@ietf.org>, and/or to the authors.

10. Acknowledgements

Thanks to Ing-jyh (Inton) Tsang for initial discussions on the need for ECN propagation in L2TP and its applicability. Thanks also to Carlos Pignataro, Tom Herbert, Ignacio Goyret, Alia Atlas, Praveen Balasubramanian, Joe Touch, Mohamed Boucadair, David Black, Jake Holland and Sri Gundavelli for helpful advice and comments. "A Comparison of IPv6-over-IPv4 Tunnel Mechanisms" [[RFC7059](#)] helped to identify a number of tunnelling protocols to include within the scope of this document.

Bob Briscoe was part-funded by the Research Council of Norway through the TimeIn project. The views expressed here are solely those of the authors.

11. References

11.1. Normative References

- [I-D.ietf-tsvwg-ecn-encap-guidelines]
Briscoe, B., Kaippallimalil, J., and P. Thaler,
"Guidelines for Adding Congestion Notification to
Protocols that Encapsulate IP", [draft-ietf-tsvwg-ecn-encap-guidelines-13](#) (work in progress), May 2019.
- [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>>.

Briscoe
18]

Expires January 9, 2020

[Page

- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", [RFC 2474](#), DOI 10.17487/RFC2474, December 1998, <<https://www.rfc-editor.org/info/rfc2474>>.
- [RFC2661] Townsley, W., Valencia, A., Rubens, A., Pall, G., Zorn, G., and B. Palter, "Layer Two Tunneling Protocol "L2TP"", [RFC 2661](#), DOI 10.17487/RFC2661, August 1999, <<https://www.rfc-editor.org/info/rfc2661>>.
- [RFC2784] Farinacci, D., Li, T., Hanks, S., Meyer, D., and P. Traina, "Generic Routing Encapsulation (GRE)", [RFC 2784](#), DOI 10.17487/RFC2784, March 2000, <<https://www.rfc-editor.org/info/rfc2784>>.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", [RFC 3168](#), DOI 10.17487/RFC3168, September 2001, <<https://www.rfc-editor.org/info/rfc3168>>.
- [RFC3931] Lau, J., Ed., Townsley, M., Ed., and I. Goyret, Ed., "Layer Two Tunneling Protocol - Version 3 (L2TPv3)", [RFC 3931](#), DOI 10.17487/RFC3931, March 2005, <<https://www.rfc-editor.org/info/rfc3931>>.
- [RFC4301] Kent, S. and K. Seo, "Security Architecture for the Internet Protocol", [RFC 4301](#), DOI 10.17487/RFC4301, December 2005, <<https://www.rfc-editor.org/info/rfc4301>>.
- [RFC4380] Huitema, C., "Teredo: Tunneling IPv6 over UDP through Network Address Translations (NATs)", [RFC 4380](#), DOI 10.17487/RFC4380, February 2006, <<https://www.rfc-editor.org/info/rfc4380>>.
- [RFC5129] Davie, B., Briscoe, B., and J. Tay, "Explicit Congestion Marking in MPLS", [RFC 5129](#), DOI 10.17487/RFC5129, January 2008, <<https://www.rfc-editor.org/info/rfc5129>>.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", [RFC 6040](#), DOI 10.17487/RFC6040, November 2010, <<https://www.rfc-editor.org/info/rfc6040>>.

11.2. Informative References

- [GTPv1] 3GPP, "GPRS Tunnelling Protocol (GTP) across the Gn and Gp interface", Technical Specification TS 29.060.

- [GTPv1-U] 3GPP, "General Packet Radio System (GPRS) Tunnelling Protocol User Plane (GTPv1-U)", Technical Specification TS 29.281.
- [GTPv2-C] 3GPP, "Evolved General Packet Radio Service (GPRS) Tunnelling Protocol for Control plane (GTPv2-C)", Technical Specification TS 29.274.
- [I-D.ietf-intarea-gue]
Herbert, T., Yong, L., and O. Zia, "Generic UDP Encapsulation", [draft-ietf-intarea-gue-07](#) (work in progress), March 2019.
- [I-D.ietf-intarea-tunnels]
Touch, J. and M. Townsley, "IP Tunnels in the Internet Architecture", [draft-ietf-intarea-tunnels-09](#) (work in progress), July 2018.
- [I-D.ietf-nvo3-geneve]
Gross, J., Ganga, I., and T. Sridhar, "Geneve: Generic Network Virtualization Encapsulation", [draft-ietf-nvo3-geneve-13](#) (work in progress), March 2019.
- [I-D.ietf-nvo3-vxlan-gpe]
Maino, F., Kreeger, L., and U. Elzur, "Generic Protocol Extension for VXLAN", [draft-ietf-nvo3-vxlan-gpe-07](#) (work in progress), April 2019.
- [I-D.ietf-sfc-nsh-ecn-support]
Eastlake, D., Briscoe, B., and A. Malis, "Explicit Congestion Notification (ECN) and Congestion Feedback Using the Network Service Header (NSH)", [draft-ietf-sfc-nsh-ecn-support-01](#) (work in progress), July 2019.
- [RFC1701] Hanks, S., Li, T., Farinacci, D., and P. Traina, "Generic Routing Encapsulation (GRE)", [RFC 1701](#), DOI 10.17487/RFC1701, October 1994, <<https://www.rfc-editor.org/info/rfc1701>>.
- [RFC2637] Hamzeh, K., Pall, G., Verthein, W., Taarud, J., Little, W., and G. Zorn, "Point-to-Point Tunneling Protocol (PPTP)", [RFC 2637](#), DOI 10.17487/RFC2637, July 1999, <<https://www.rfc-editor.org/info/rfc2637>>.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", [RFC 2983](#), DOI 10.17487/RFC2983, October 2000, <<https://www.rfc-editor.org/info/rfc2983>>.

Briscoe
20]

Expires January 9, 2020

[Page

- [RFC3260] Grossman, D., "New Terminology and Clarifications for Diffserv", [RFC 3260](#), DOI 10.17487/RFC3260, April 2002, <<https://www.rfc-editor.org/info/rfc3260>>.
- [RFC3308] Calhoun, P., Luo, W., McPherson, D., and K. Peirce, "Layer Two Tunneling Protocol (L2TP) Differentiated Services Extension", [RFC 3308](#), DOI 10.17487/RFC3308, November 2002, <<https://www.rfc-editor.org/info/rfc3308>>.
- [RFC5415] Calhoun, P., Ed., Montemurro, M., Ed., and D. Stanley, Ed., "Control And Provisioning of Wireless Access Points (CAPWAP) Protocol Specification", [RFC 5415](#), DOI 10.17487/RFC5415, March 2009, <<https://www.rfc-editor.org/info/rfc5415>>.
- [RFC5845] Muhanna, A., Khalil, M., Gundavelli, S., and K. Leung, "Generic Routing Encapsulation (GRE) Key Option for Proxy Mobile IPv6", [RFC 5845](#), DOI 10.17487/RFC5845, June 2010, <<https://www.rfc-editor.org/info/rfc5845>>.
- [RFC5944] Perkins, C., Ed., "IP Mobility Support for IPv4, Revised", [RFC 5944](#), DOI 10.17487/RFC5944, November 2010, <<https://www.rfc-editor.org/info/rfc5944>>.
- [RFC6275] Perkins, C., Ed., Johnson, D., and J. Arkko, "Mobility Support in IPv6", [RFC 6275](#), DOI 10.17487/RFC6275, July 2011, <<https://www.rfc-editor.org/info/rfc6275>>.
- [RFC6830] Farinacci, D., Fuller, V., Meyer, D., and D. Lewis, "The Locator/ID Separation Protocol (LISP)", [RFC 6830](#), DOI 10.17487/RFC6830, January 2013, <<https://www.rfc-editor.org/info/rfc6830>>.
- [RFC7059] Steffann, S., van Beijnum, I., and R. van Rein, "A Comparison of IPv6-over-IPv4 Tunnel Mechanisms", [RFC 7059](#), DOI 10.17487/RFC7059, November 2013, <<https://www.rfc-editor.org/info/rfc7059>>.
- [RFC7296] Kaufman, C., Hoffman, P., Nir, Y., Eronen, P., and T. Kivinen, "Internet Key Exchange Protocol Version 2 (IKEv2)", STD 79, [RFC 7296](#), DOI 10.17487/RFC7296, October 2014, <<https://www.rfc-editor.org/info/rfc7296>>.

Briscoe
21]

Expires January 9, 2020

[Page

- [RFC7348] Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger, L., Sridhar, T., Bursell, M., and C. Wright, "Virtual eXtensible Local Area Network (VXLAN): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks", [RFC 7348](#), DOI 10.17487/RFC7348, August 2014, <<https://www.rfc-editor.org/info/rfc7348>>.
- [RFC7450] Bumgardner, G., "Automatic Multicast Tunneling", [RFC 7450](#), DOI 10.17487/RFC7450, February 2015, <<https://www.rfc-editor.org/info/rfc7450>>.
- [RFC7637] Garg, P., Ed. and Y. Wang, Ed., "NVGRE: Network Virtualization Using Generic Routing Encapsulation", [RFC 7637](#), DOI 10.17487/RFC7637, September 2015, <<https://www.rfc-editor.org/info/rfc7637>>.
- [RFC7665] Halpern, J., Ed. and C. Pignataro, Ed., "Service Function Chaining (SFC) Architecture", [RFC 7665](#), DOI 10.17487/RFC7665, October 2015, <<https://www.rfc-editor.org/info/rfc7665>>.
- [RFC8085] Eggert, L., Fairhurst, G., and G. Shepherd, "UDP Usage Guidelines", [BCP 145](#), [RFC 8085](#), DOI 10.17487/RFC8085, March 2017, <<https://www.rfc-editor.org/info/rfc8085>>.
- [RFC8087] Fairhurst, G. and M. Welzl, "The Benefits of Using Explicit Congestion Notification (ECN)", [RFC 8087](#), DOI 10.17487/RFC8087, March 2017, <<https://www.rfc-editor.org/info/rfc8087>>.
- [RFC8159] Konstantynowicz, M., Ed., Heron, G., Ed., Schatzmayr, R., and W. Henderickx, "Keyed IPv6 Tunnel", [RFC 8159](#), DOI 10.17487/RFC8159, May 2017, <<https://www.rfc-editor.org/info/rfc8159>>.
- [RFC8229] Pauly, T., Touati, S., and R. Mantha, "TCP Encapsulation of IKE and IPsec Packets", [RFC 8229](#), DOI 10.17487/RFC8229, August 2017, <<https://www.rfc-editor.org/info/rfc8229>>.
- [RFC8300] Quinn, P., Ed., Elzur, U., Ed., and C. Pignataro, Ed., "Network Service Header (NSH)", [RFC 8300](#), DOI 10.17487/RFC8300, January 2018, <<https://www.rfc-editor.org/info/rfc8300>>.

Briscoe
22]

Expires January 9, 2020

[Page

Internet-Draft
2019

ECN over IP-shim-(L2)-IP Tunnels

July

Author's Address

Bob Briscoe
Independent
UK

E-Mail: ietf@bobbriscoe.net

URI: <http://bobbriscoe.net/>

Briscoe
23]

Expires January 9, 2020

[Page