

Internet Area WG
Internet Draft
Intended status: Informational
Updates: [4459](#)
Expires: January 2017

J. Touch
USC/ISI
M. Townsley
Cisco
July 6, 2016

IP Tunnels in the Internet Architecture
draft-ietf-intarea-tunnels-03.txt

Status of this Memo

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

This document may contain material from IETF Documents or IETF Contributions published or made publicly available before November 10, 2008. The person(s) controlling the copyright in some of this material may not have granted the IETF Trust the right to allow modifications of such material outside the IETF Standards Process. Without obtaining an adequate license from the person(s) controlling the copyright in such materials, this document may not be modified outside the IETF Standards Process, and derivative works of it may not be created outside the IETF Standards Process, except to format it for publication as an RFC or to translate it into languages other than English.

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/1id-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 January 6, 2017.

Copyright Notice

Copyright (c) 2016 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](http://trustee.ietf.org/license-info) 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.

Abstract

This document discusses the role of IP tunnels in the Internet architecture, in which IP datagrams are carried as payloads in non-link layer protocols. It explains their relationship to existing protocol layers and the challenges in supporting IP tunneling based on the equivalence of tunnels to links.

Table of Contents

1.	Introduction.....	3
2.	Conventions used in this document.....	6
2.1.	Key Words.....	6
2.2.	Terminology.....	6
3.	The Tunnel Model.....	9
3.1.	What is a tunnel?.....	10
3.2.	View from the Outside.....	11
3.3.	View from the Inside.....	12
3.4.	Location of the Ingress and Egress.....	12
3.5.	Implications of This Model.....	13
3.6.	Fragmentation.....	14
3.6.1.	Outer Fragmentation.....	14
3.6.2.	Inner Fragmentation.....	15
3.6.3.	The necessity of Outer Fragmentation.....	16
4.	IP Tunnel Requirements.....	16
4.1.	Minimum MTU Considerations.....	17
4.2.	Fragmentation.....	18
4.3.	MTU discovery.....	21
4.4.	IP ID exhaustion.....	22
4.5.	Hop Count.....	23
4.6.	Signaling.....	24

4.7.	Relationship of Header Fields.....	26
4.8.	Congestion.....	27
4.9.	Checksums.....	27
4.10.	Numbering.....	27
4.11.	Multicast.....	28
4.12.	Multipoint.....	28
4.13.	NAT / Load Balancing.....	29
4.14.	Recursive tunnels.....	29
5.	Observations (implications).....	29
5.1.	Tunnel protocol designers.....	29
5.2.	Tunnel implementers.....	30
5.3.	Tunnel operators.....	30
5.4.	Diagnostics.....	30
5.5.	For existing standards.....	31
5.5.1.	Generic UDP Encapsulation (GUE - IP in UDP in IP)...	31
5.5.2.	Generic Packet Tunneling in IPv6.....	31
5.5.3.	Geneve (NV03).....	32
5.5.4.	GRE (IP in GRE in IP).....	33
5.5.5.	IP in IP / mobile IP.....	33
5.5.6.	IPsec tunnel mode (IP in IPsec in IP).....	35
5.5.7.	L2TP.....	36
5.5.8.	L2VPN.....	36
5.5.9.	L3VPN.....	36
5.5.10.	LISP.....	36
5.5.11.	MPLS.....	37
5.5.12.	PWE.....	37
5.5.13.	SEAL/AERO.....	37
5.5.14.	TRILL.....	37
5.5.15.	RTG DT encapsulations.....	38
5.6.	For future standards.....	38
6.	Security Considerations.....	39
7.	IANA Considerations.....	40
8.	References.....	40
8.1.	Normative References.....	40
8.2.	Informative References.....	40
9.	Acknowledgments.....	44
APPENDIX A:	Fragmentation efficiency.....	45
A.1.	Selecting fragment sizes.....	45
A.2.	Packing.....	46

1. Introduction

The Internet is loosely based on the ISO seven layer stack, in which data units traverse the stack by being wrapped inside data units one layer down. A tunnel is a mechanism for transmitting data units between endpoints by wrapping them as data units of the same or higher layers, e.g., IP in IP (Figure 1) or IP in UDP (Figure 2).

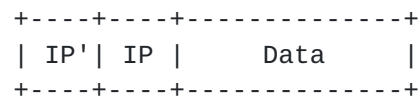


Figure 1 IP inside IP

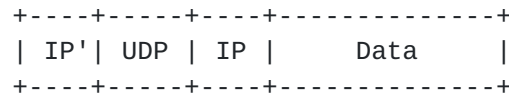


Figure 2 IP in UDP in IP in Ethernet

This document focuses on tunnels that transit IP packets, i.e., in which an IP packet is the payload of another protocol. Tunnels provide a virtual link that can help decouple the network topology seen by transiting packets from the underlying physical network [To98][RFC2473]. Tunnels were critical in the development of multicast because not all routers were capable of processing multicast packets [Er94]. Tunnels allowed multicast packets to transit between multicast-capable routers over paths that did not support multicast. Similar techniques have been used to support other protocols, such as IPv6 [RFC2460].

Use of tunnels is common in the Internet. The word "tunnel" occurs in over 100 RFCs, and is supported within numerous protocols, including:

- o IP in IP / mobile IP - IPv4 in IPv4 tunnels [RFC2003][RFC2473][RFC5944]
- o IP in IPv6 - IPv6 or IPv4 in IPv6 [RFC2473]
- o IPsec - includes a tunnel mode to enable encryption or authentication of the an entire IP datagram [RFC4301]
- o Generic Router Encapsulation (GRE) - a shim layer for tunneling any network layer in any other network layer, IP in GRE in IP [RFC2784][RFC7588][RFC7676]
- o Generic UDP Encapsulation (GUE) - IP in UDP (in IP)[He15]
- o Automatic Multicast Tunneling (AMT) - IP in UDP for multicast [RFC7450]
- o L2TP - PPP over IP, to extend a subscriber's DSL/FTTH connection from an access line provider to an ISP [RFC3931]

- o L2VPNs - provides a link topology different from that provided by physical links [[RFC4664](#)]
- o L3VPNs - provides a network topology different from that provided by ISPs [[RFC4176](#)]
- o LISP - reduces routing table load within an enclave of routers at the expense of more complex ingress encapsulation tables [[RFC6830](#)]
- o MPLS - IP over a circuit-like path in which identifiers are rewritten on each hop, often used for traffic provisioning [[RFC3031](#)]
- o NV03 - data center network sharing (which includes use of GUE, above) [[RFC7364](#)]
- o PWE3 - emulates wire-like services over packet-switched services [[RFC3985](#)]
- o SEAL/AERO -IP in IP tunneling with an additional shim header designed to overcome the limitations of [RFC2003](#) [[RFC5320](#)][Te16]
- o TRILL - enables L3 routing (typically IS-IS) in an enclave of Ethernet bridges [[RFC5556](#)][RFC6325]

The variety of tunnel mechanisms raises the question of the role of tunnels in the Internet architecture and the potential need for these mechanisms to have similar and predictable behavior. In particular, the ways in which packet sizes (i.e., Maximum Transmission Unit or MTU) mismatch and error signals (e.g., ICMP) are handled may benefit from a coordinated approach.

Regardless of the layer in which encapsulation occurs, tunnels emulate a link. The only difference is that a link operates over a physical communication channel, whereas a tunnel operates over software protocol layers. Because tunnels are links, they are subject to the same issues as any link, e.g., MTU discovery, signaling, and the potential utility of native support for broadcast and multicast [[RFC2460](#)][RFC3819]. They have advantages over native links, being potentially easier to reconfigure and control.

The first attempt to use large-scale tunnels transit multicast across the Internet in 1988 lead to tunnel collapse. At the time, tunnels were not implemented as encapsulation-based virtual links, but rather as loose source routes on un-encapsulated IP datagrams [[RFC1075](#)]. Using encapsulation tunnels instead avoided that collapse [[Er94](#)] and eventually to AMT [[RFC7450](#)].

The remainder of this document describes the general principles of IP tunneling and discusses the key considerations in the design of a protocol that tunnels IP datagrams. It derives its conclusions from the equivalence of tunnels and links. Note that all considerations are in the context of existing standards and requirements.

2. Conventions used in this document

2.1. Key Words

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](#)].

2.2. Terminology

This document uses the following terminology. These definitions are given in the most general terms, but will be used primarily to discuss IP tunnels in this document. They are presented in order from most fundamental to those derived on earlier definitions:

- o Messages: variable length data labeled with globally-unique endpoint IDs, also known as a datagram for IP messages [[RFC791](#)].
- o Network node (node): a device that can act as an endpoint or forwarder. For datagrams (IP messages), these are hosts or gateways/routers, respectively.
- o Endpoint or host: a node that sources or sinks messages labeled from/to its IDs, typically known as a host for both IP and higher-layer protocol messages [[RFC1122](#)].
- o Forwarder: a node that relays messages using destination IDs and local context, also known as a gateway or router for IP messages [[RFC1812](#)]. Note that most forwarders also act as endpoints when they source or sink messages.
- o Source (sender): the node that generates a message.
- o Destination (receiver): the node that consumes a message.
- o Link: a device (or medium) that transfers messages between nodes, i.e., by which a message can traverse between nodes without being processed by a forwarder. Note that the notion of forwarder is relative to the layer at which message processing is considered [[To16](#)].

- o Link interface (sometimes known as a network interface): a location on a link co-located with a node where messages depart onto that link or arrive from that link.
- o Path: a sequence of one or more links or tunnels over which a message can traverse between nodes (hosts or forwarders), which may or may not involve being processed by a forwarder.
- o Tunnel: a protocol mechanism that transits messages using encapsulation to allow a path to appear as a single link. Note that a protocol can be used to tunnel itself (IP over IP) and that this includes the conventional layering of the ISO stack (i.e., by this definition, Ethernet is a tunnel for IP). A tunnel can be considered a virtual link.
- o Ingress: the virtual link interface of a tunnel which receives messages within a node, encapsulates them according to the tunnel protocol, and transmits them into the tunnel. This is the tunnel equivalent of the outgoing (departing) network interface of a link. Note that the ingress virtual link interface and traffic source node can be co-located.
- o Egress: a virtual link interface that receives messages that have finished transiting a tunnel and presents them to a node. This is the tunnel equivalent of the incoming (arriving) network interface of a link. The egress decapsulates messages for further transit to the destination. Note that the egress virtual link interface and traffic destination node can be co-located.
- o Tunnel transit packet (TTP): the packet arriving at a node connected to a tunnel that enters the ingress and exits the egress, i.e., the packet carried over the tunnel. This is sometimes known as the "tunneled packet", i.e., the packet carried over the tunnel. This is the tunnel equivalent of a network layer packet as it would traverse a link.
- o Tunnel link packet (TLP): packets that traverse from ingress to egress, in which resides all or part of a tunnel transit packet. This is sometimes known as the "tunnel packet", i.e., the packet of the tunnel itself. This is the tunnel equivalent of a link layer packet as it would traverse a link.

- o Link MTU (LMTU): the largest message that can transit a link. It typically does not include link-layer information, e.g., link layer headers or trailers, i.e., it refers to the message that the link can carry rather than the message as it appears on the link. This is thus the largest network layer packet (including network layer headers, e.g., IP datagram) that can transit a link. Note that this need not be the native size of messages on the link, i.e., the link may internally fragment and reassemble messages. For IPv4, the smallest LMTU is 68 bytes [[RFC791](#)], and for IPv6 the smallest LMTU is 1280 bytes [[RFC2460](#)].
- o Path MTU (PMTU): the largest message that can transit a path. Typically, this is the minimum of the link MTUs of the links of the path, and represents the largest network layer message (including network layer headers) that can transit a path. Note that this is not the largest network packet that can be sent between a source and destination; this is the largest network network packet that can be sent without requiring reassembly at the network layer of the destination.
- o Reassembly MTU (RMTU): the largest message that can be reassembled by a destination, which is not directly related to the link or path MTU. Sometimes also referred to as "receiver MTU". For IPv4, this is 576 bytes [[RFC793](#)] and for IPv6 it is 1500 bytes [[RFC2460](#)]; note that in both cases, the size refers to the message transferred at the network layer, which includes the network layer headers.
- o Tunnel MTU (TMTU): the largest message that can transit a tunnel, i.e., this is the tunnel equivalent of a link MTU. Typically, this is limited by the egress reassembly MTU. Note that this value may have no relation to the path MTU between the tunnel ingress and egress.
- o Tunnel internal MTU (TIMTU): the largest message that a tunnel egress can emit into a tunnel without requiring further fragmentation to reach the tunnel egress. This the path MTU between the ingress and egress.
- o Egress reassembly MTU (ERMTU): the largest message that can be reassembled by an egress. This is the size of the RMTU of a tunnel minus the encapsulation overhead of that tunnel. Sometimes also referred to as the "egress MTU".

3. The Tunnel Model

A network architecture is an abstract description of a distributed communications system, its components and their relationships, the requisite properties of those components and the emergent properties of the system that result [To03]. Such descriptions can help explain behavior, as when the OSI seven-layer model is used as a teaching example [Zi80]. Architectures describe capabilities - and, just as importantly, constraints.

A network can be defined as a system of endpoints and relays interconnected by communication paths, abstracting away issues of naming in order to focus on message forwarding. To the extent that the Internet has a single, coherent interpretation, its architecture is defined by its core protocols (IP [RFC791], TCP [RFC793], UDP [RFC768]) and messages, hosts, routers, and links [C188][To03], as shown in Figure 3:

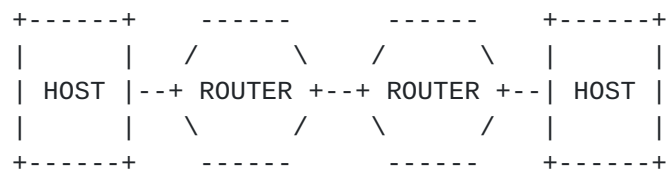


Figure 3 Basic Internet architecture

As a network architecture, the Internet is a system of hosts and routers interconnected by links that exchange messages when possible. "When possible" defines the Internet's "best effort" principle. The limited role of routers and links represents the End-to-End Principle [Sa84] and longest-prefix match enables hierarchical forwarding.

Although the definitions of host, router, and link seem absolute, they are often relative as viewed within the context of one OSI layer, each of which can be considered a distinct network architecture. An Internet gateway is a Layer 3 router when it transits IP datagrams but it acts as a Layer 2 host as it sources or sinks Layer 2 messages on attached links to accomplish this transit capability. In this way, a single node (Internet gateway) behaves as different components (router, host) at different layers.

Even though a single node may have multiple roles - even concurrently - at a given layer, each role is typically static and determined by context. An Internet gateway always acts as a Layer 2 host and that behavior does not depend on where the gateway is viewed from within Layer 2. In the context of a single layer, a node's behavior is modeled as a single component from all viewpoints in that layer.

3.1. What is a tunnel?

A tunnel can be modeled as a link in another network [To98][To01][To03]. In Figure 4, a source host (Hsrc) and destination host (Hdst) communicating over a network M in which two routers (Ra and Rd) are connected by a tunnel. Keep in mind that it is possible that both network N and network M can both be components of the Internet, i.e., there may be regular traffic as well as tunneled traffic over any of the routers shown.

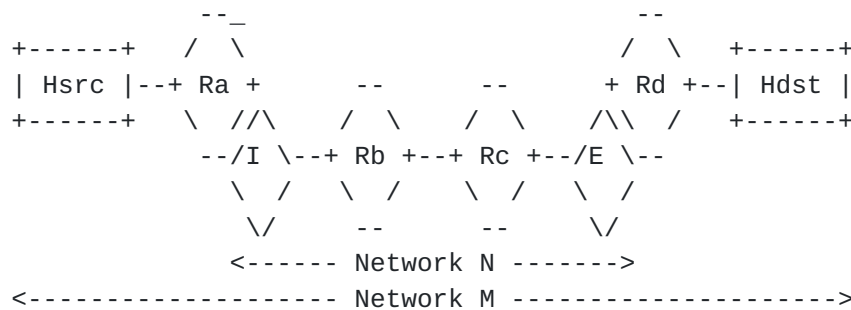


Figure 4 The big picture

The tunnel consists of two elements (ingress I, egress E), that lie along a path connected by a (possibly different) network N. Regardless of how the ingress and egress are connected, the tunnel serves as a link to the nodes it connects (here, Ra and Rd).

IP packets arriving at the ingress are encapsulated to traverse network N. We call these packets "tunnel transit packets" (TTPs) because they will now transit the tunnel inside one or more "tunnel link packets" (TLPs). TLPs use the source address of the ingress and the destination address of the egress - using whatever address is appropriate to the Layer at which the ingress and egress operate (Layer 2, Layer 3, Layer 4, etc.). The egress decapsulates those messages, which then continue on network M as if emerging from a link. To tunnel transit packets, and to the routers the tunnel connects (Ra and Rd), the tunnel acts as a link and the ingress and egress act as network interfaces to that link.

The model of each component (ingress, egress) and the entire system (tunnel) depends on the layer from which you view the tunnel. From the perspective of the outermost hosts (Hsrc and Hdst), the tunnel appears as a link between two routers (Ra and Rd). For routers along the tunnel (e.g., Rb and Rc), the ingress and egress appear as the endpoint hosts and Hsrc and Hdst are invisible.

When the tunnel network (N) is implemented using the same protocol as the endpoint network (M), the picture looks flatter (Figure 5), as if it were running over a single network. However, note that this appearance is incorrect - nothing has changed. From the perspective of the endpoints, Rb and Rc and network N don't exist and aren't visible, and from the perspective of the tunnel, network M doesn't exist. The fact that network N and M use the same protocol, and may traverse the same links is irrelevant.

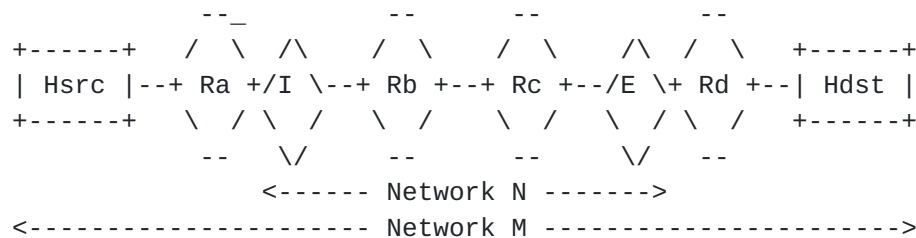


Figure 5 IP in IP network picture

3.2. View from the Outside

From outside the tunnel, to network M, the entire tunnel acts as a link (Figure 6). It may be numbered or unnumbered and the addresses associated with the ingress and egress are irrelevant from outside.



Figure 6 Tunnels as viewed from the outside

A tunnel is effectively invisible to the network in which it resides, except that it behaves exactly as a link. Consequently [RFC3819] requirements for links supporting IP also apply to tunnels.

E.g., the IP datagram hop count (IPv4 Time-to-Live [RFC791] and IPv6 Hop Limit [RFC2460]) are decremented when traversing a router, not by traversing a link - or thus a tunnel. Tunnels have a tunnel MTU - the largest datagram that can transit, just as links have a corresponding link MTU. A link MTU may not reflect the native link message sizes (ATM AAL5 48 byte messages support a 9KB MTU) and the same is true for a tunnel.

3.3. View from the Inside

Within network N, i.e., from inside the tunnel itself, the ingress is a source of tunnel link packets and the egress is a sink - both are hosts on network N (Figure 7). Consequently [\[RFC1122\]](#) Internet host requirements apply to ingress and egress nodes when Network N uses IP (and thus the ingress/egress use IP encapsulation).

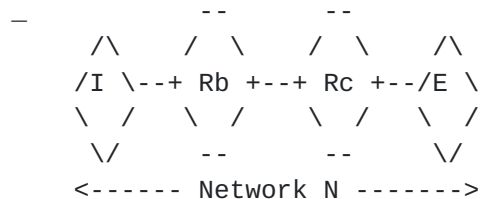


Figure 7 Tunnels, as viewed from within the tunnel

Viewed from within the tunnel, the outer network (M) doesn't exist. Tunnel link packets can be fragmented by the source (ingress) and reassembled at the destination (egress), just as at any endpoint. The path between ingress and egress may have a path MTU but the endpoints can exchange messages as large as can be reassembled at the destination (egress), i.e., an egress MTU. Information about the network - i.e., regarding MTU sizes, network reachability, etc. - are relayed from the destination (egress) and intermediate routers back to the source (ingress), without regard for the external network (M).

3.4. Location of the Ingress and Egress

The ingress and egress are endpoints of the tunnel and the tunnel is a link. The ingress and egress are thus link endpoints at the network nodes the tunnel interconnects. Such link endpoints are typically described as "network interfaces".

Tunnel interfaces may be physical or virtual. The interface may be implemented inside the node where the tunnel attaches, e.g., inside a host or router. The interface may also be implemented as a "bump in the wire" (BITW), somewhere along a link between the two nodes the link interconnects. IP in IP tunnels are often implemented as interfaces, where IPsec tunnels are sometimes implemented as BITW. These implementation variations determine only whether information available at the link endpoints (ingress/egress) can be easily shared with the connected network nodes.

3.5. Implications of This Model

This approach highlights a few key features of a tunnel as a network architecture construct:

- o To the tunnel transit packets (TTPs), tunnels turn a network (Layer 3) path into a (Layer 2) link
- o To nodes the tunnel traverses, the tunnel ingress and egress act as hosts that source and sink tunnel link packets (TLPs)

The consequences of these features are as follow:

- o Like a link, a tunnel has an MTU defined by the reassembly MTU of the receiving interface (egress).
- o Like any other link, the MTU inside a tunnel are not relevant to the transited traffic. There is no mechanism or protocol by which they are measured or confirmed.
- o Path MTU discovery in the network layer (i.e., outer network M) has no direct relation to the MTU of the hops within the link layer of the links (or thus tunnels) that connect its components.
- o Hops remain defined as the number of routers encountered on a path or the time spent at a router [[RFC1812](#)]. Hops are not decremented solely by the transit of a link, e.g., a packet with a hop count of zero should successfully transit a link (and thus a tunnel) that connects two hosts. Routers, not links, alter hopcounts.
- o The addresses of a tunnel ingress and egress correspond to link layer addresses to the tunnel transit packet and outer network M. Like point-to-point links, point-to-point tunnels can be unnumbered in the network in which they reside (even though they must have addresses in the network they transit).
- o Like network interfaces, the ingress and egress are never a direct source of ICMP messages but may provide information to their attached host or router to generate those ICMP messages.
- o Like network interfaces and links, two nodes may be connected by any combination of tunnels and links, including multiple tunnels. As with multiple links, existing routing determines which traffic uses each link or tunnel.

These observations make it much easier to determine what a tunnel must do to transit IP packets, notably it must satisfy all

requirements expected of a link [RFC1122][RFC3819]. The consequence of these observations are that tunnels are no different from links, except only that a link has a physical instantiation.

3.6. Fragmentation

There are two places where fragmentation can occur in a tunnel, called Outer Fragmentation and Inner Fragmentation. This document assumes that only Outer Fragmentation is viable because it is the only approach that works for IPv4 datagrams with DF=1 and for IPv6.

3.6.1. Outer Fragmentation

The simplest case is Outer Fragmentation, as shown in Figure 8. The bottom of the figure shows the network topology, where packets start at the source, enter the tunnel at the encapsulator, exit the tunnel at the decapsulator, and arrive finally at the destination. The packet traffic is shown above the topology, where the end-to-end packets are shown at the top. The packets are composed of an inner header (iH) and inner data (iD); the term "inner") is relative to the tunnel, as will become apparent. When the packet (iH,iD) arrives at the encapsulator, it is placed inside the tunnel packet structure, here shown as adding just an outer header, oH, in step (a).

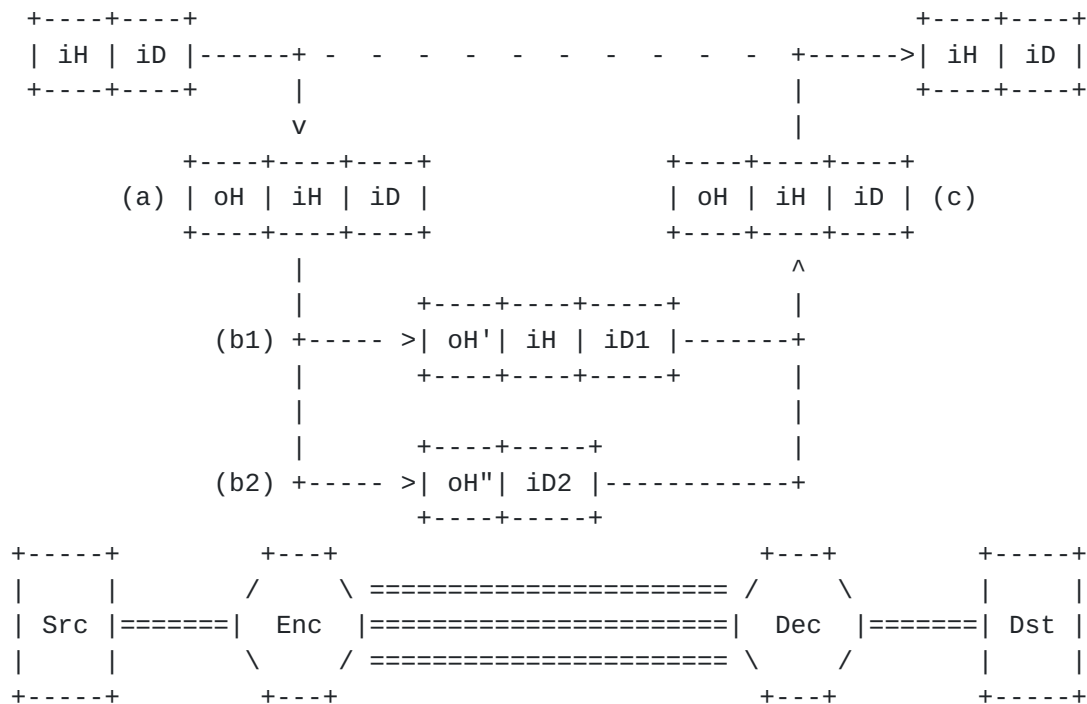


Figure 8 Fragmentation of the outer packet

When the encapsulated packet exceeds the tunnel MTU, the packet needs to be fragmented. In this case we fragment the packet at the outer header, with the fragments shown as (b1) and (b2). Note that the outer header indicates fragmentation (as ' and "), the inner header occurs only in the first fragment, and the inner data is broken across the two packets. These fragments are reassembled at the encapsulator in step (c), and the resulting packet is decapsulated and sent on to the destination.

Outer fragmentation isolates Source and Destination from tunnel encapsulation duties. This can be considered a benefit in clean, layered network design, but also may result in complex decapsulator design, especially where tunnels aggregate large amounts of traffic, such as IP ID overload (see Sec. 4.4). Outer fragmentation is valid for any tunnel encapsulation protocol that supports fragmentation (e.g., IPv4 or IPv6), where the tunnel endpoints act as the host endpoints of that protocol.

Along the tunnel, the inner header is contained only in the first fragment, which can interfere with mechanisms that 'peek' into lower layer headers, e.g., as for ICMP, as discussed in Sec. 4.6.

3.6.2. Inner Fragmentation

Inner Fragmentation distributes the impact of tunneling across both the decapsulator and destination, and is shown in Figure 9; this can be especially important when the tunnel aggregates large amounts of traffic. However, this mechanism is thus valid only when the original source packets can be fragmented on-path, e.g., as in IPv4 datagrams with DF=0.

Again, the network topology is shown at the bottom of the figure, and the original packets show at the top. Packets arrive at the encapsulator, and are fragmented there based on the inner header into (a1) and (a2). The fragments arrive at the decapsulator, which removes the outer header and forwards the resulting fragments on to the destination. The destination is then responsible for reassembling the fragments into the original packet.

Along the tunnel, the inner headers are copied into each fragment, and so are available to mechanisms that 'peek' into headers (e.g., ICMP, as discussed in Sec. 4.6). Because fragmentation happens on the inner header, the impact of IP ID is reduced.

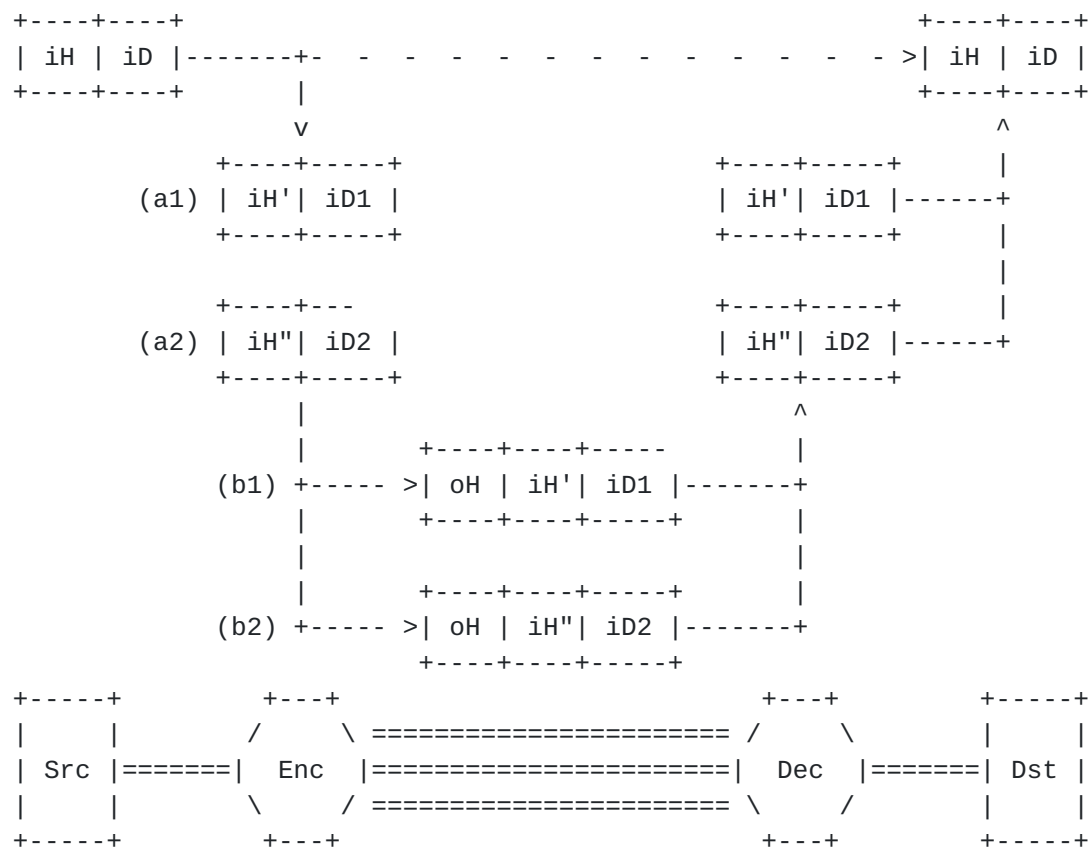


Figure 9 Fragmentation of the inner packet

3.6.3. The necessity of Outer Fragmentation

Fragmentation is critical tunnels that support TTP packets for protocols with minimum MTU requirements, while operating over tunnel paths using protocols with minimum MTU requirements. Depending on the amount of space used by encapsulation, these two minimums will ultimately interfere, and the TTP will need to be fragmented to both support a TTP minimum MTU while traversing tunnels with their own TLP minimum MTUs.

Outer Fragmentation is the only solution that supports all IPv4 and IPv6 traffic, because inner fragmentation is allowed only for IPv4 datagrams with DF=0. As a result, the remainder of this document assumes Outer Fragmentation.

4. IP Tunnel Requirements

The requirements of an IP tunnel are defined by the requirements of an IP link because both transit IP packets. A tunnel thus must transit the IP minimum MTU, i.e., 68 bytes for IPv4 [RFC793] and 1280

bytes for IPv6 [[RFC2460](#)] and a tunnel must support address resolution when there is more than one egress.

The requirements of the tunnel ingress and egress are defined by the network over which they exchange messages (tunnel link packets). For IP-over-IP, this means that the ingress **MUST NOT** exceed the IPv4 Identification (fragment) field uniqueness requirements [[RFC6864](#)].

These requirements remain even though tunnels have some unique issues, including the need for additional space for encapsulation headers and the potential for tunnel MTU variation.

4.1. Minimum MTU Considerations

There are a variety of values of minimum MTU to consider, both in a conventional network and in a tunnel as a link in that network. These are indicated in Figure 10, an annotated variant of Figure 4.

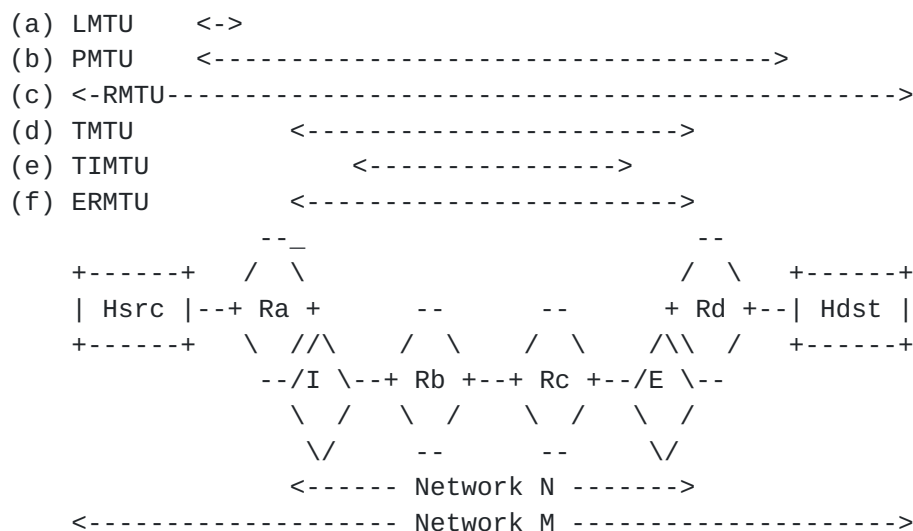


Figure 10 The variety of MTU values

Consider the following example values. For IPv6, the minimum LMTU (a) is 1280 bytes, which is also the minimum PMTU (b). The minimum RMTU (c) is 1500 bytes, which is also the minimum MTU for endpoint-to-endpoint communication. This means that IPv6 already assumes that endpoint-to-endpoint communication may require source fragmentation to transit IPv6-compatible links, even without considering tunnels.

The TMTU (d) is the tunnel equivalent of a LMTU, and thus also needs to be 1280 bytes for IPv6. Assuming the links of a tunnel traverse IPv6 hops (e.g., I to Rb, Rb to Rc, and Rc to E), the TIMTU (e) is equivalent to the PMTU between I and E, which is 1280 - encaps (where

"encaps" is the tunnel encapsulation overhead). This value is insufficient to satisfy the requirement of an IPv6 link (which must transit at least 1280 bytes unfragmented), but this is not a problem. The TMTU (d) is not limited by TIMTU (e), but by ERMTU (f), the tunnel equivalent of RMTU (c). For a tunnel using IPv6 over IPv6, the ERMTU is the RMTU of the underlying network N minus space for encapsulation, i.e., 1500 - encaps bytes, and the tunnel is viable as long as ERMTU \geq 1280. Even though the tunnel will ultimately transit ERMTU - encaps byte messages between the ingress and egress, each hop within the tunnel transits only TIMTU - encaps byte messages. The difference between TIMTU and ERMTU is the reason why the tunnel ingresses need to support fragmentation and tunnel egresses need to support reassembly. The high cost of fragmentation and reassembly is why it is useful for applications to avoid sending messages too close to the PMTU, even the PMTU at their own layer.

[4.2. Fragmentation](#)

A tunnel interacts with fragmentation in two different ways. As a link in network M, its messages might be fragmented before they reach the tunnel - i.e., at the TTP layer either during source fragmentation (if generated at the same node as the ingress interface) or forwarding fragmentation (for IPv4 DF=0 datagrams). In addition, messages traversing the tunnel may require fragmentation by the ingress - i.e., source fragmentation at the TLP layer by the ingress. These two fragmentation operations are no more related than are conventional IP fragmentation and ATM segmentation and reassembly; one occurs at the network layer, the other at the (virtual) link layer.

As with any link layer, a tunnel MTU (TMTU) is defined as the largest message that can transit the tunnel. For a tunnel, this is the egress reassembly MTU (ERMTU), which is the reassembly MTU (RMTU) of the egress interface minus the space needed for the tunnel encapsulation headers. This value must also satisfy the requirements of the IP packets that the tunnel transits.

Note that many of the issues with tunnel fragmentation and MTU handling were discussed in [\[RFC4459\]](#), but that document described a variety of alternatives as if they were independent. This document explains the combined approach that is necessary.

Like any other link, an IPv4 tunnel must transit 68 byte packets without requiring source fragmentation [\[RFC791\]](#)[\[RFC1122\]](#) and an IPv6 tunnel must transit 1280 byte packets without requiring source fragmentation [\[RFC2460\]](#). The tunnel MTU interacts with routers or hosts it connects the same way as would a link MTU. In the following

pseudocode, TTPsize is the size of the tunnel transit packet (TTP), and ERMTU is the reassembly MTU of the egress. As with any link, the link MTU (LMTU) is defined not by the native path of the link (or, for a tunnel, the path MTU of encapsulated packets inside the tunnel) but by the egress reassembly capability. This is because the ICMP "packet too big" message indicates failure of a link to transit a packet, not a preference for a size that matches that inside the mechanism of the link. There is no ICMP message for "larger than I'd like, but I can still transit it".

These rules apply at the host/router where the tunnel is attached, i.e., at the network layer of the TTP (we assume that all tunnels, including multipoint tunnels, have a single, uniform TMTU). These are basic source fragmentation rules (or transit refragmentation for IPv4 DF=0 datagrams), and have no relation to the tunnel itself other than to consider the TMTU as the effective LMTU of the next hop:

```
if (TTP > TMTU) then
  if (TTP can be fragmented, e.g., IPv4 DF=0) then
    split TTP into fragments of TMTU size
    and send each fragment to the tunnel ingress
  else
    drop TTP and send ICMP "too big" to TTP source
  endif
else
  send TTP to the tunnel ingress
endif
```

These rules apply at the tunnel ingress, in its role as host on the tunnel path, i.e., as source fragmentation of TTP messages (we assume that all tunnels, even multipoint tunnels, have a single, uniform TIMTU), where "encaps" is the encapsulation overhead:

```
if (TTP <= (TIMTU + encaps)) then
  encapsulate the TTP and process as if arriving at the node
else
  if ((TIMTU + encaps) < TTP <= (ERMTU - encaps)) then
    fragment TTP into TIMTU chunks
    encapsulate each chunk and process as if arriving at the node
  else
    {never happens; host/router already dropped by now}
  endif
endif
```


There is one path above that never occurs - i.e., a network interface should never receive a message larger than its MTU, and a tunnel should thus never receive a message larger than its (ERMTU - encaps) limit. A router attempting to process such a message would generate an ICMP error (packet too big, fragmentation needed) and the packet would already have been dropped before entering into this algorithm.

As an example, consider IPv4 over IPv6 or IPv6 over IPv6 tunneling, where IPv6 encapsulation adds a 40 byte fixed header plus IPv6 options (i.e., IPv6 header extensions) of total size TOptSz. From [\[RFC2460\]](#) it follows that the TMTU must be at least 1280 bytes and the ERMTU must be at least 1500 - (40 + TOptSz) bytes. The TIMTU must be a minimum of 1280 - (40 + TOptSz) bytes. Considering these minimum values, the previous algorithm becomes:

```
if (TTP <= (1240 - TOptSz)) then
    encapsulate the TTP and process as if arriving at the node
else
    if ((1240 - TOptSz) < TTP <= (1460 - TOptSz)) then
        fragment TTP into (1240 - TOptSz) chunks
        encapsulate each chunk and process as if arriving at the node
    else
        {never happens; host/router already dropped by now}
    endif
endif
```

This tunnel supports IPv6 transit only if TOptSize is smaller than 180 bytes, and supports IPv4 transit if TOptSize is smaller than 884 bytes. IPv6 TTPs of 1280 bytes may be guaranteed transit the outer network (M) without needing fragmentation there but they may require ongoing fragmentation and reassembly if the TMTU is not at least 1320 bytes.

When using IP directly over IP, the minimum ERMTU for IPv4 is 576 bytes and for IPv6 is 1500 bytes. This means that tunnels of IPv4-over-IPv4, IPv4-over-IPv6, and IPv6-over-IPv6 are possible without additional requirements, but this may involve ingress fragmentation and egress reassembly. IPv6 cannot be tunneled directly over IPv4 without additional requirements, notably that the ERMTU is at least 1280 bytes. Fragmentation and reassembly cannot be avoided for IPv6-over-IPv6 without similar requirements.

When ongoing ingress fragmentation and egress reassembly would be prohibitive or costly, larger MTUs can be supported by design and confirmed either out-of-band (by design) or in-band (e.g., using PLPMTUD [\[RFC4821\]](#), as done in SEAL [\[RFC5320\]](#) and AERO [\[Te16\]](#)).

Alternately, an ingress can encapsulate packets that fit and shut down once fragmentation is needed, but it must not continue to forward smaller packets while dropping larger packets that are still within required limits.

4.3. MTU discovery

MTU discovery enables a network path to support a larger PMTU than it can assume from the minimum requirements of protocol over which it operates. A tunnel has two different LMTU-like values: TMTU and the TIMTU.

There is temptation to optimize tunnel traversal so that packets are not fragmented between ingress and egress, i.e., to attempt tune the network PMTU to the TIMTU rather than the TMTU, to avoid ingress fragmentation. This is hazardous for many reasons:

- o The tunnel is capable of transiting packets as large as the ERMTU, which is always at least as large as the TIMTU and typically is larger.
- o ICMP has only one type of error message regarding large packets - "too big", i.e., too large to transit. There is no optimization message of "bigger than I'd like, but I can deal with if needed".
- o IP tunnels often involve some level of recursion, i.e., encapsulation over itself [[RFC4459](#)].

Recursive tunneling occurs whenever a protocol ends up encapsulated in itself. This happens directly, as when IPv4 is encapsulated in IPv4, or indirectly, as when IP is encapsulated in UDP which then is a payload inside IP. It can involve many layers of encapsulation because a tunnel provider isn't always aware of whether the packets it transits are already tunneled.

Recursion is impossible when the tunnel transit packets are limited to that of the native size of the TIMTU. Arriving tunnel transit packets have a minimum supported size (1280 for IPv6) and the tunnel PMTU has the same requirement; there would be no room for the additional encapsulation headers. The result would be an IPv6 tunnel that cannot satisfy IPv6 transit requirements.

It is more appropriate to require the tunnel to satisfy IP transit requirements and enforce that requirement at design time or during operation (the latter using PLPMTUD [[RFC4821](#)]). Conventional path MTU discovery (PMTUD) relies on existing endpoint ICMP processing of explicit negative feedback from routers along the path via "message

to big" ICMP packets in the reverse direction of the tunnel [[RFC1191](#)]. This technique is susceptible to the "black hole" phenomenon, in which the ICMP messages never return to the source due to policy-based filtering [[RFC2923](#)]. PLPMTUD requires a separate, direct control channel from the egress to the ingress that provides positive feedback; the direct channel is not blocked by policy filters and the positive feedback ensures fail-safe operation if feedback messages are lost [[RFC4821](#)].

4.4. IP ID exhaustion

In IPv4, the IP Identification (ID) field is a 16-bit value that is unique for every packet for a given source address, destination address, and protocol, such that it does not repeat within the Maximum Segment Lifetime (MSL) [[RFC791](#)][[RFC1122](#)]. Although the ID field was originally intended for fragmentation and reassembly, it can also be used to detect and discard duplicate packets, e.g., at congested routers (see Sec. 3.2.1.5 of [[RFC1122](#)]). For this reason, and because IPv4 packets can be fragmented anywhere along a path, all packets between a source and destination of a given protocol must have unique ID values over a period of an MSL, which is typically interpreted as two minutes (120 seconds). These requirements have recently been somewhat relaxed in recognition of the primary use of this field for reassembly and the need to handle only fragment misordering at the receiver [[RFC6864](#)].

The uniqueness of the IP ID is a known problem for high speed nodes, because it limits the speed of a single protocol between two endpoints [[RFC4963](#)]. Although this suggests that the uniqueness of the IP ID is moot, tunnels exacerbate this condition. A tunnel often aggregates traffic from a number of different source and destination addresses, of different protocols, and encapsulates them in a header with the same ingress and egress addresses, all using a single encapsulation protocol. The result is one of the following:

1. The IP ID rules are enforced, and the tunnel throughput is severely limited.
2. The IP ID rules are enforced, and the tunnel consumes large numbers of ingress/egress IP addresses solely to ensure ID uniqueness.
3. The IP ID rules are ignored.

The last case is the most obvious solution, because it corresponds to how endpoints currently behave. Fortunately, fragmentation is somewhat rare in the current Internet at large, but it can be common

along a tunnel. Fragments that repeat the IP ID risk being reassembled incorrectly, especially when fragments are reordered or lost. Reassembly errors are not always detected by other protocol layers (see Sec. 4.9), and even when detected they can result in excessive overall packet loss and can waste bandwidth between the egress and ultimate packet destination.

4.5. Hop Count

This section considers the selection of the value of the hop count of the tunnel link header, as well as the potential impact on the tunnel transit header. The former is affected by the number of hops within the tunnel. The latter determines whether the tunnel has visible effect on the transit packet.

In general, the Internet hop count field is used to detect and avoid forwarding loops that cannot be corrected without a synchronized reboot. The IPv4 Time-to-Live (TTL) and IPv6 Hop Limit field each serve this purpose [[RFC791](#)][RFC2460].

The IPv4 TTL field was originally intended to indicate packet expiration time, measured in seconds. A router is required to decrement the TTL by at least one or the number of seconds the packet is delayed, whichever is larger [[RFC1812](#)]. Packets are rarely held that long, and so the field has come to represent the count of the number of routers traversed. IPv6 makes this meaning more explicit.

These hop count fields represent the number of network forwarding elements traversed by an IP datagram. An IP datagram with a hop count of zero can traverse a link between two hosts because it never visits a router (where it would need to be decremented and would have been dropped).

An IP datagram traversing a tunnel thus need not have its hopcount modified, i.e., the tunnel transit header need not be affected. A zero hop count datagram should be able to traverse a tunnel as easily as it traverses a link. A router MAY be configured to decrement packets traversing a particular link (and thus a tunnel), which may be useful in emulating a path as if it had traversed one or more routers, but this is strictly optional. The ability of the outer network and tunnel network to avoid indefinitely looping packets does not rely on the hop counts of the tunnel traversal packet and tunnel link packet being related in any way at all.

The hop count field is also used by several protocols to determine whether endpoints are "local", i.e., connected to the same subnet (link-local discovery and related protocols [[RFC4861](#)]). A tunnel is a

way to make a remote address appear directly-connected, so it makes sense that the other ends of the tunnel appear local and that such link-local protocols operate over tunnels unless configured explicitly otherwise. When the interfaces of a tunnel are numbered, these can be interpreted the same way as if they were on the same link subnet.

4.6. Signaling

In the current Internet architecture, signaling goes upstream, either from routers along a path or from the destination, back toward the source. Such signals are typically contained in ICMP messages, but can involve other protocols such as RSVP, transport protocol signals (e.g., TCP RSTs), or multicast control or transport protocols.

A tunnel behaves like a link and acts like a link interface at the nodes where it is attached. As such, it can provide information that enhances IP signaling (e.g., ICMP), but itself does not directly generate ICMP messages.

For tunnels, this means that there are two separate signaling paths. The outer network M nodes can each signal the source of the tunnel transit packets, Hsrc (Figure 11). Inside the tunnel, the inner network N nodes can signal the source of the tunnel link packets, the ingress I (Figure 12).

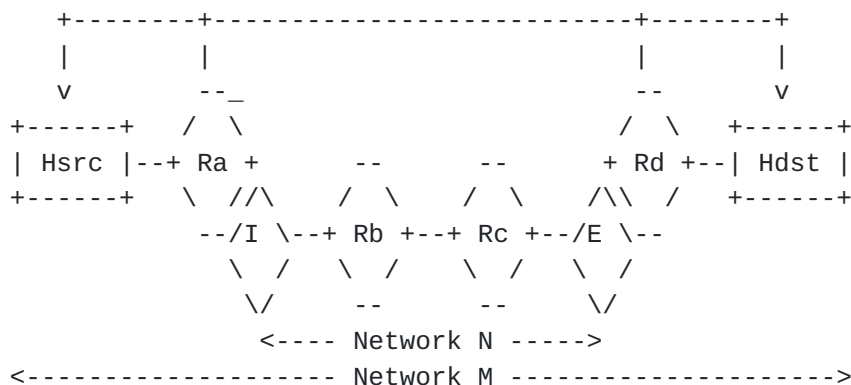


Figure 11 Signals outside the tunnel

4.7. Relationship of Header Fields

Some tunnel specifications attempt to relate the fields of the tunnel transit packet and tunnel link packet, i.e., the packet arriving at the ingress and the encapsulation header. These two headers are effectively independent and there is no utility in requiring their contents to be related.

In specific, the encapsulation header source and destination addresses are network endpoints in the tunnel network N, but have no meaning in the outer network M, even when the tunneled packet traverses the same network. The addresses are effectively independent, and the tunnel endpoint addresses are link addresses to the tunnel transit packet.

Because the tunneled packet uses source and destination addresses with a separate meaning, it is inappropriate to copy or reuse the IPv4 Identification or IPv6 Fragment ID fields of the tunnel transit packet. These fields need to be generated based on the context of the encapsulation header, not the tunnel transit header.

Similarly, the DF field need not be copied from the tunnel transit packet to the encapsulation header of the tunnel link packet (presuming both are IPv4). Path MTU discovery inside the tunnel does not directly correspond to path MTU discovery outside the tunnel, i.e., inside the tunnel it would update the TIMTU used for outer fragmentation at the ingress, but has no effect on the TMTU reported to the device where the ingress is attached as a network interface.

The same is true for most other fields. When a field value is generated in the encapsulation header, its meaning should be derived from what is desired in the context of the tunnel as a link. When feedback is received from these fields, they should be presented to the tunnel ingress and egress as if they were network interfaces. The behavior of the node where these interfaces attach should be identical to that of a conventional link.

There are exceptions to this rule that are explicitly intended to relay signals from inside the tunnel to outside the tunnel. The primary example is ECN [[RFC6040](#)], which copies the ECN bits from the tunnel transit header to the tunnel link header during encapsulation at the ingress and modifies the tunnel transit header at egress based on a combination of the bits of the two headers. This is intended to allow congestion notification within the tunnel to be interpreted as if it were on the direct path. Other examples may involve the DSCP flags. In both cases, it is assumed that the intent of copying values on encapsulation and merging values on decapsulation has the effect

of allowing the tunnel to act as if it participates in the same type of network as outside the tunnel (network M).

4.8. Congestion

In general, tunnels carrying IP traffic need not react directly to congestion any more than would any other link layer [[RFC5405](#)]. IP traffic is not generally expected to be congestion reactive.

[text from David Black on ECN relaying?]

4.9. Checksums

IP traffic transiting a tunnel needs to expect a similar level of error detection and correction as it would expect from any other link. In the case of IPv4, there are no such expectations, which is partly why it includes a header checksum [[RFC791](#)].

IPv6 omitted the header checksum because it already expects most link errors to be detected and dropped by the link layer and because it also assumes transport protection [[RFC2460](#)]. When transiting IPv6 over IPv6, the tunnel fails to provide the expected error detection. This is why IPv6 is often tunneled over layers that include separate protection, such as GRE [[RFC2784](#)].

The fragmentation created by the tunnel ingress can increase the need for stronger error detection and correction, especially at the tunnel egress to avoid reassembly errors. The Internet checksum is known to be susceptible to reassembly errors that could be common [[RFC4963](#)], and should not be relied upon for this purpose. This is why SEAL and AERO include a separate checksum [[RFC5320](#)][Te16]. This requirement can be undermined when using UDP as a tunnel with no UDP checksum (as per [[RFC6935](#)][RFC6936]) when fragmentation occurs because the egress has no checksum with which to validate reassembly. For this reason, it is safe to use UDP with a zero checksum for atomic (non-fragmented, non-fragmentable) tunnel link packets only; when used on fragments, whether generated at the ingress or en-route inside the tunnel, omission of such a checksum can result in reassembly errors that can cause additional work (capacity, forwarding processing, receiver processing) downstream of the egress.

4.10. Numbering

Tunnel ingresses and egresses have addresses associated with the encapsulation protocol. These addresses are the source and destination (respectively) of the encapsulated packet while traversing the tunnel network.

Tunnels may or may not have addresses in the network whose traffic they transit (e.g., network M in Figure 4). In some cases, the tunnel is an unnumbered interface to a point-to-point virtual link. When the tunnel has multiple egresses, tunnel interfaces require separate addresses in network M.

To see the effect of tunnel interface addresses, consider traffic sourced at router Ra in Figure 4. Even before being encapsulated by the ingress, that traffic needs a source IP network address that belongs to the router. One option is to use an address associated with one of the other interfaces of the router [[RFC1122](#)]. Another option is to assign a number to the tunnel interface itself. Regardless of which address is used, the resulting IP packet is then encapsulated by the tunnel ingress using the ingress address as a separate operation.

4.11. Multicast

[To be addressed]

Note that PMTU for multicast is difficult. PIM carries an option that may help in the Population Count Extensions to PIM [[RFC6807](#)].

IMO, again, this is no different than any other multicast link.

4.12. Multipoint

Multipoint tunnels are tunnels with more than two ingress/egress endpoints. Just as tunnels emulate links, multipoint tunnels emulate multipoint links.

Multipoint links require a support for egress determination, just as multipoint links do. This function is typically supported by ARP [[RFC826](#)] or ARP emulation (e.g., LAN Emulation, known as LANE [[RFC2225](#)]) for multipoint links. For multipoint tunnels, a similar mechanism is required for the same purpose - to determine the egress address for proper ingress encapsulation.

All multipoint systems - tunnels and links - might support different MTUs between each ingress/egress (or link entrance/exit) pair. In most cases, it is simpler to assume a uniform MTU throughout the multipoint system, e.g., the minimum MTU supported across all ingress/egress pairs. This applies to both the ERMTU and TIMETU (the latter as used only by the ingress).

A multipoint tunnel **MUST** have support for broadcast and multicast, in exactly the same way as this is already required for multipoint links

[[RFC3819](#)]. Both modes can be supported either by a native mechanism inside the tunnel or by emulation using serial replication at the tunnel ingress, in the same way that links may provide the same support either natively (e.g., via promiscuous or automatic replication in the link itself) or network interface emulation (e.g., as for non-broadcast multiaccess networks, i.e., NBMA).

[4.13.](#) NAT / Load Balancing

[To be addressed]

Talk about ECMP / LAG here

[4.14.](#) Recursive tunnels

[IS THIS REDUNDANT?]

The rules described in this document already support tunnels over tunnels, sometimes known as "recursive" tunnels, in which IP is transited over IP either directly or via intermediate encapsulation (IP-UDP-IP).

There are known hazards to recursive tunneling, notably that the independence of the tunnel transit header and tunnel link header hop counts can result in a tunneling loop. Such looping can be avoided when using direct encapsulation (IP in IP) by use of a header option to track the encapsulation count and to limit that count [[RFC2473](#)]. This looping cannot be avoided when other protocols are used for tunneling, e.g., IP in UDP in IP, because the encapsulation count may not be visible where the recursion occurs.

[5.](#) Observations (implications)

[Leave this as a shopping list for now]

[5.1.](#) Tunnel protocol designers

Recursive tunneling + minimum MTU = frag/reassembly is inevitable, at least to be able to split/join two fragments

Account for egress MTU/path MTU differences.

Include a stronger checksum.

Ensure the egress MTU is always larger than the path MTU.

Ensure that the egress reassembly can keep up with line rate OR design PLPMTUD into the tunneling protocol.

5.2. Tunnel implementers

Detect when the egress MTU is exceeded.

Detect when the egress MTU drops below the required minimum and shut down the tunnel if that happens - configuring the tunnel down and issuing a hard error may be the only way to detect this anomaly, and it's sufficiently important that the tunnel SHOULD be disabled. This is always better than blindly assuming the tunnel has been deployed correctly, i.e., that the solution has been engineered.

Do NOT decrement the TTL as part of being a tunnel. It's always already OK for a router to decrement the TTL based on different next-hop routers, but TTL is a property of a router not a link.

5.3. Tunnel operators

Keep the difference between "enforced by operators" vs. "enforced by active protocol mechanism" in mind. It's fine to assume something the tunnel cannot or does not test, as long as you KNOW you can assume it. When the assumption is wrong, it will NOT be signaled by the tunnel. Do NOT decrement the TTL as part of being a tunnel. It's always already OK for a router to decrement the TTL based on different next-hop routers, but TTL is a property of a router not a link.

Do NOT decrement the TTL as part of being a tunnel. It's always already OK for a router to decrement the TTL based on different next-hop routers, but TTL is a property of a router not a link.

>>>> PLPMTUD can give incorrect information during ECMP or LAG

5.4. Diagnostics

Some current implementations include diagnostics to support monitoring the impact of tunneling, especially the impact on fragmentation and reassembly resources, the status of path MTU discovery, etc.

>> Because a tunnel ingress/egress is a network interface, it SHOULD have similar resources as any other network interface. This includes resources for packet processing as well as monitoring.

5.5. For existing standards

5.5.1. Generic UDP Encapsulation (GUE - IP in UDP in IP)

[He15]

Consistent with this doc:

Inconsistent with this doc:

Imports [RFC4459](#)

Appears to allow both pre and post-encapsulation fragmentation

Recommendations:

Should not encourage pre-encaps fragmentation

See recommendations for [RFC4459](#)

5.5.2. Generic Packet Tunneling in IPv6

[RFC2473]

Consistent with this doc:

Considers the endpoints of the tunnel as virtual interfaces.

Considers the tunnel a virtual link.

Requires source fragmentation at the ingress and reassembly at the egress.

Includes a recursion limit to prevent unlimited re-encapsulation.

Sets tunnel transit header hop limit independently.

Sends ICMPs back at the ingress based on the arriving tunnel transit packet and its relation to the tunnel MTU (though it uses the incorrect value of the tunnel MTU; see below).

Allows for ingress relaying of internal tunnel errors (but see below; it does not discuss retaining state about these).

Inconsistent with this doc:

Decrements the tunnel transit header by 1, i.e., incorrectly assuming that tunnel endpoints occur at routers only and that the tunnel, rather than the router, is responsible for this decrement.

This doc goes to pains to describe the decapsulation process as if it were distinct from conventional protocol processing by the receiver (when it should not be).

Copies traffic class from tunnel link to tunnel transit header (as one variant).

Treats the tunnel MTU as the tunnel path MTU, rather than the tunnel egress MTU.

Incorrectly fragments IPv4 DF=0 tunnel transit packets that arrive larger than the tunnel MTU at the IPv6 layer; the relationship between IPv4 and the tunnel is more complex (as noted in this doc).

Fails to retain state from the tunnel based on ingress receiving ICMP messages from inside the tunnel, e.g., such as might cause future tunnel transit packets arriving at the ingress to be discarded with an ICMP error response rather than allowing them to proceed into the tunnel.

Recommendation:

This doc should update 2473 for TTL decrement, tunnel MTU, and fragmentation. Other issues are less critical.

5.5.3. Geneve (NV03)

[RFC7364] info, [[Gr16](#)] stds - ISSUE US AS BCP; Gr16 should follow

Consistent with this doc:

Generation of the link header fields is not discussed and presumed independent of transit packet.

Reportedly treats an ingress/egress as applying to multiple tunnels, rather than considering them logically independent for each tunnel. This appears to confuse implementation aggregation with architecture.

Reportedly treats tunnels as supporting traffic for multiple virtual networks, rather than considering them logically independent. This appears to confuse implementation aggregation with architecture.

Inconsistent with this doc:

Tries to match transit to tunnel path MTU rather than egress MTU.

Recommendation:

Gr16 should be updated to follow us

5.5.4. GRE (IP in GRE in IP)

IPv4 [[RFC2784](#)] stds, [[RFC7588](#)] info, [[RFC7676](#)] stds - NO CHANGES

Consistent with this doc:

Does not address link header generation.

Non-default behavior allows fragmentation of link packet to match tunnel path MTU up to the limit of the egress MTU.

Default behavior sets link DF independently.

Shuts the tunnel down if the tunnel path MTU isn't ≥ 1280 .

Inconsistent with this doc:

Based on tunnel path MTU, not egress MTU.

Claims that the tunnel (GRE) mechanism is responsible for generating ICMP error messages.

Default behavior fragments transit packet (where possible) based on tunnel path MTU (it should fragment based on egress MTU).

Default behavior does not support the minimum MTU of IPv6 when run over IPv6.

Non-default behavior allows copying DF for IPv4 in IPv4.

Recommendations:

No changes - existing docs largely describe legacy deployment.

5.5.5. IP in IP / mobile IP

IPv4 [[RFC2003](#)] stds, [[RFC4459](#)] info:

Consistent with this doc:

Generate link ID independently

Generate link DF independently when transit DF=0

Generate ECN/update ECN based on sharing info [[RFC6040](#)]

Set link TTL to transit to egress only (independently)

Do not decrement TTL on entry except when part of forwarding

Do not decrement TTL on exit except when part of forwarding

Options not copied, but used as a hint to desired services.

Generally treat tunnel as a link, e.g., for link-local.

Inconsistent with this doc

Set link DF when transit DF=1 (won't work unless I-E runs PLPMTUD)

Drop at egress if transit TTL=0 (wrong TTL for host-host tunnels)

Drop when transit source is router's IP (prevents tun from router)

Drop when transit source matches egress (prevents tun to router)

Use tunnel ICMPs to generate upper ICMPs, copying context (ICMPs are now coming from inside a link!); these should be handled by setting errors as a "network interface" and letting the attached host/router figure out what to send.

Using tunnel MTU discovery to tune the transit packet to the tunnel path MTU rather than egress MTU.

Recommendations:

IMO, ought to update 2003! (no "update" to informational), esp. regarding TTL issues, transit source drop issues, and tunnel MTU.

IPv6 [[RFC2473](#)] std:

Consistent with this doc:

Doesn't discuss lots of header fields, but implies they're set independently.

Sets link TTL independently.

Inconsistent with this doc:

Tunnel issues ICMP PTBs.

ICMP PTB issued if larger than 1280 - header, rather than egress reassembly MTU.

Fragments IPv6 over IPv6 fragments only if transit is ≤ 1280 (i.e., forces all tunnels to have a max MTU of 1280).

Fragments IPv4 over IPv6 fragments only if IPv4 DF=0 (misinterpreting the "can fragment the IPv4 packet" as permission to fragment at the IPv6 link header)

Considers encapsulation a forwarding operation and decrements the transit TTL.

Recommendation:

Should UPDATE 2473; tunnel should not issue PTBs (router should), issue them correctly, fragment correctly, and not TTL decrement.

[5.5.6. IPsec tunnel mode \(IP in IPsec in IP\)](#)

[RFC4301] std

Consistent with this doc:

Most of the rules, except as noted below.

Inconsistent with this doc:

Writes its own header copying rules (Sec 5.1.2), rather than referring to existing standards, but that makes sense for security reasons.

Uses policy to set, clear, or copy DF (policy isn't the issue)

Intertwines tunneling with forwarding rather than presenting the tunnel as a network interface; this can be corrected by using IPsec transport mode with an IP-in-IP tunnel [[RFC3884](#)].

Recommendations:

None.

5.5.7. L2TP

[RFC3931] std

Consistent with this doc:

Does not address most link headers, which are thus independent.

Inconsistent with this doc:

Manages tunnel access based on tunnel path MTU, instead of egress MTU.

Refers to [RFC2473](#) (IPv6 in IPv6), which is inconsistent with this doc as noted above.

Recommendations:

Should update to use correct tunnel MTU.

5.5.8. L2VPN

[RFC4664]

Consistent with this doc:

Inconsistent with this doc:

Recommendations:

5.5.9. L3VPN

[RFC4176]

Consistent with this doc:

Inconsistent with this doc:

Recommendations:

5.5.10. LISP

[RFC6830]

Consistent with this doc:

Inconsistent with this doc:

Recommendations:

[5.5.11](#). MPLS

[RFC3031]

Consistent with this doc:

Inconsistent with this doc:

Recommendations:

[5.5.12](#). PWE

[RFC3985]

Consistent with this doc:

Inconsistent with this doc:

Recommendations:

[5.5.13](#). SEAL/AERO

[[RFC5320](#)][Te16]

Consistent with this doc:

Inconsistent with this doc:

Recommendations:

[5.5.14](#). TRILL

[[RFC5556](#)][RFC6325]

Consistent with this doc:

Puts IP in Ethernet, so most of the issues don't come up.

Ethernet doesn't have TTL or fragment.

Rbridge (trill) TTL header is independent of transit packet.

Inconsistent with this doc:

None.

Recommendations:

None.

5.5.15. RTG DT encapsulations

[[No16](#)], refers to NV03 and other encapsulations

Includes info on tables for multipoint tunnels, additional info for headers, etc.

Consistent with this doc:

Inconsistent with this doc:

Assumes MTU can be managed to avoid fragmentation. This is impossible as long as any one layer is used recursively and that layer includes a mandatory minimum MTU. A "trust but verify" policy is better than assuming engineered MTU deployment is sufficient.

Relies on ICMP PTB to correct for tunnel path MTU issues.

Allows encaps protocols to not support fragmentation.

Recommendations:

That doc should refer to this regarding general tunneling issues, including fragmentation, tunnel MTU, and TTL, including the "trust but verify" issue for engineered MTU deployment.

All encaps protocols for IP over IP (eventually) MUST support fragm.

5.6. For future standards

Larger IPv4 MTU (2K? or just 2x path MTU?) for reassembly

Always include frag support for at least two frags; do NOT try to deprecate fragmentation.

Limit encapsulation option use/space.

Augment ICMP to have two separate messages: PTB vs P-bigger-than-optimal

Include MTU as part of BGP as a hint - SB

Hazards of multi-MTU [draft-van-beijnum-multi-mtu-04](#)

6. Security Considerations

Tunnels may introduce vulnerabilities or add to the potential for receiver overload and thus DOS attacks. These issues are primarily related to the fact that a tunnel is a link that traverses a network path and to fragmentation and reassembly. ICMP signal translation introduces a new security issue and must be done with care. ICMP generation at the router or host attached to a tunnel is already covered by existing requirements (e.g., should be throttled).

Tunnels traverse multiple hops of a network path from ingress to egress. Traffic along such tunnels may be susceptible to on-path and off-path attacks, including fragment injection, reassembly buffer overload, and ICMP attacks. Some of these attacks may not be as visible to the endpoints of the architecture into which tunnels are deployed and these attacks may thus be more difficult to detect.

Fragmentation at routers or hosts attached to tunnels may place an undue burden on receivers where traffic is not sufficiently diffuse, because tunnels may induce source fragmentation at hosts and path fragmentation (for IPv4 DF=0) more for tunnels than for other links. Care should be taken to avoid this situation, notably by ensuring that tunnel MTUs are not significantly different from other link MTUs.

Tunnel ingresses emitting IP datagrams MUST obey all existing IP requirements, such as the uniqueness of the IP ID field. Failure to either limit encapsulation traffic, or use additional ingress/egress IP addresses, can result in high speed traffic fragments being incorrectly reassembled.

Tunnels are susceptible to attacks at both the inner and outer network layers. The tunnel ingress/egress endpoints appear as network interfaces in the outer network, and are as susceptible as any other network interface. This includes vulnerability to fragmentation reassembly overload, traffic overload, and spoofed ICMPs that misreport the state of those interfaces. Similarly, the ingress/egress appear as hosts to the path traversed by the tunnel, and thus are as susceptible as any other host to attacks as well.

[management?]

[Access control?]

describe relationship to [\[RFC6169\]](#) - JT (as per INTAREA meeting notes, don't cover Teredo-specific issues in [RFC6169](#), but include generic issues here)

7. IANA Considerations

This document has no IANA considerations.

The RFC Editor should remove this section prior to publication.

8. References

8.1. Normative References

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

8.2. Informative References

- [Cl88] Clark, D., "The design philosophy of the DARPA internet protocols," Proc. Sigcomm 1988, p.106-114, 1988.
- [Er94] Eriksson, H., "MBone: The Multicast Backbone," Communications of the ACM, Aug. 1994, pp.54-60.
- [Gr16] Gross, J., et al., "Geneve: Generic Network Virtualization Encapsulation," [draft-ietf-nvo3-geneve-01](#), Jan. 2016.
- [He15] Herbert, T., L. Yong, O. Zia, "Generic UDP Encapsulation," [draft-ietf-nvo3-gue-04](#), Jul. 2016.
- [No16] Nordmark, E. (Ed.), A. Tian, J. Gross, J. Hudson, L. Kreeger, P. Garg, P. Thaler, T. Herbert, "Encapsulation Considerations," [draft-ietf-rtgwg-dt-encap-01](#), Mar. 2016.
- [RFC768] Postel, J., "User Datagram Protocol," [RFC 768](#), Aug. 1980
- [RFC791] Postel, J., "Internet Protocol," [RFC 791](#) / STD 5, September 1981.
- [RFC793] Postel, J., "Transmission Control Protocol," [RFC 793](#), Sept. 1981.
- [RFC826] Plummer, D., "An Ethernet Address Resolution Protocol -- or -- Converting Network Protocol Addresses to 48.bit Ethernet Address for Transmission on Ethernet Hardware," [RFC 826](#), Nov. 1982.

- [RFC1075] Waitzman, D., C. Partridge, S. Deering, "Distance Vector Multicast Routing Protocol," [RFC 1075](#), Nov. 1988.
- [RFC1122] Braden, R., Ed., "Requirements for Internet Hosts - Communication Layers," [RFC 1122](#) / STD 3, October 1989.
- [RFC1191] Mogul, J., S. Deering, "Path MTU discovery," [RFC 1191](#), November 1990.
- [RFC1812] Baker, F., "Requirements for IP Version 4 Routers," [RFC 1812](#), June 1995.
- [RFC2003] Perkins, C., "IP Encapsulation within IP," [RFC 2003](#), October 1996.
- [RFC2225] Laubach, M., J. Halpern, "Classical IP and ARP over ATM," [RFC 2225](#), Apr. 1998.
- [RFC2460] Deering, S., R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification," [RFC 2460](#), Dec. 1998.
- [RFC2473] Conta, A., "Generic Packet Tunneling in IPv6 Specification," [RFC 2473](#), Dec. 1998.
- [RFC2784] Farinacci, D., T. Li, S. Hanks, D. Meyer, P. Traina, "Generic Routing Encapsulation (GRE)", [RFC 2784](#), March 2000.
- [RFC2923] Lahey, K., "TCP Problems with Path MTU Discovery," [RFC 2923](#), September 2000.
- [RFC2473] Conta, A., S. Deering, "Generic Packet Tunneling in IPv6 Specification," [RFC 2473](#), Dec. 1998.
- [RFC3031] Rosen, E., A. Viswanathan, R. Callon, "Multiprotocol Label Switching Architecture", [RFC 3031](#), January 2001.
- [RFC3819] Karn, P., Ed., C. Bormann, G. Fairhurst, D. Grossman, R. Ludwig, J. Mahdavi, G. Montenegro, J. Touch, L. Wood, "Advice for Internet Subnetwork Designers," [RFC 3819](#) / [BCP 89](#), July 2004.
- [RFC3884] Touch, J., L. Eggert, Y. Wang, "Use of IPsec Transport Mode for Dynamic Routing," [RFC 3884](#), September 2004.

- [RFC3931] Lau, J., Ed., M. Townsley, Ed., I. Goyret, Ed., "Layer Two Tunneling Protocol - Version 3 (L2TPv3)," [RFC 3931](#), March 2005.
- [RFC3985] Bryant, S., P. Pate (Eds.), "Pseudo Wire Emulation Edge-to-Edge (PWE3) Architecture", [RFC 3985](#), March 2005.
- [RFC4176] El Mghazli, Y., Ed., T. Nadeau, M. Boucadair, K. Chan, A. Gonguet, "Framework for Layer 3 Virtual Private Networks (L3VPN) Operations and Management," [RFC 4176](#), October 2005.
- [RFC4301] Kent, S., and K. Seo, "Security Architecture for the Internet Protocol," [RFC 4301](#), December 2005.
- [RFC4459] Savola, P., "MTU and Fragmentation Issues with In-the-Network Tunneling," [RFC 4459](#), April 2006.
- [RFC4664] Andersson, L., Ed., E. Rosen, Ed., "Framework for Layer 2 Virtual Private Networks (L2VPNs)," [RFC 4664](#), September 2006.
- [RFC4821] Mathis, M., J. Heffner, "Packetization Layer Path MTU Discovery," [RFC 4821](#), March 2007.
- [RFC4861] Narten, T., E. Nordmark, W. Simpson, H. Soliman, "Neighbor Discovery for IP version 6 (IPv6)," [RFC 4861](#), Sept. 2007.
- [RFC4963] Heffner, J., M. Mathis, B. Chandler, "IPv4 Reassembly Errors at High Data Rates," [RFC 4963](#), July 2007.
- [RFC5320] Templin, F., Ed., "The Subnetwork Encapsulation and Adaptation Layer (SEAL)," [RFC 5320](#), Feb. 2010.
- [RFC5405] Eggert, L., G. Fairhurst, "Unicast UDP Usage Guidelines for Application Designers," [RFC 5405](#), Nov. 2008.
- [RFC5556] Touch, J., R. Perlman, "Transparently Interconnecting Lots of Links (TRILL): Problem and Applicability Statement," [RFC 5556](#), May 2009.
- [RFC5944] Perkins, C., Ed., "IP Mobility Support for IPv4, Revised" [RFC 5944](#), Nov. 2010.
- [RFC6040] Briscoe, B., "Tunneling of Explicit Congestion Notification," [RFC 6040](#), Nov. 2010.

- [RFC6169] Krishnan, S., D. Thaler, J. Hoagland, "Security Concerns With IP Tunneling," [RFC 6169](#), Apr. 2011.
- [RFC6325] Perlman, R., D. Eastlake, D. Dutt, S. Gai, A. Ghanwani, "Routing Bridges (RBridges): Base Protocol Specification," [RFC 6325](#), July 2011.
- [RFC6807] Farinacci, D., G. Shepherd, S. Venaas, Y. Cai, "Population Count Extensions to Protocol Independent Multicast (PIM)," [RFC 6807](#), Dec. 2012.
- [RFC6830] Farinacci, D., V. Fuller, D. Meyer, D. Lewis, "The Locator/ID Separation Protocol," [RFC 6830](#), Jan. 2013.
- [RFC6864] Touch, J., "Updated Specification of the IPv4 ID Field," Proposed Standard, [RFC 6864](#), Feb. 2013.
- [RFC6935] Eubanks, M., P. Chimento, M. Westerlund, "IPv6 and UDP Checksums for Tunneled Packets," [RFC 6935](#), Apr. 2013.
- [RFC6936] Fairhurst, G., M. Westerlund, "Applicability Statement for the Use of IPv6 UDP Datagrams with Zero Checksums," [RFC 6936](#), Apr. 2013.
- [RFC7364] Narten, T., Gray, E., Black, D., Fang, L., Kreeger, L., M. Napierala, "Problem Statement: Overlays for Network Virtualization", [RFC 7364](#), Oct. 2014.
- [RFC7450] Bumgardner, G., "Automatic Multicast Tunneling," [RFC 7450](#), Feb. 2015.
- [RFC7588] Bonica, R., C. Pignataro, J. Touch, "A Widely-Deployed Solution to the Generic Routing Encapsulation Fragmentation Problem," [RFC 7588](#), July 2015.
- [RFC7676] Pignataro, C., R. Bonica, S. Krishnan, "IPv6 Support for Generic Routing Encapsulation (GRE)," [RFC 7676](#), Oct 2015.
- [Sa84] Saltzer, J., D. Reed, D. Clark, "End-to-end arguments in system design," ACM Trans. on Computing Systems, Nov. 1984.
- [Te16] Templin, F., "Asymmetric Extended Route Optimization," [draft-templin-aerolink-67](#), Jun. 2016.
- [To01] Touch, J., "Dynamic Internet Overlay Deployment and Management Using the X-Bone," Computer Networks, July 2001, pp. 117-135.

- [To03] Touch, J., Y. Wang, L. Eggert, G. Finn, "Virtual Internet Architecture," USC/ISI Tech. Report 570, Aug. 2003.
- [To16] Touch, J., "Middleboxes Models Compatible with the Internet," USC/ISI Tech. Report <TBD>, July 2016.
- [To98] Touch, J., S. Hotz, "The X-Bone," Proc. Globecom Third Global Internet Mini-Conference, Nov. 1998.
- [Zi80] Zimmermann, H., "OSI Reference Model - The ISO Model of Architecture for Open Systems Interconnection," IEEE Trans. on Comm., Apr. 1980.

9. Acknowledgments

This document originated as the result of numerous discussions among the authors, Jari Arkko, Stuart Bryant, Lars Eggert, Ted Faber, Gorry Fairhurst, Dino Farinacci, Matt Mathis, and Fred Templin. It benefitted substantially from detailed feedback from Toerless Eckert, Vincent Roca, and Lucy Yong, as well as other members of the Internet Area Working Group.

This document was prepared using 2-Word-v2.0.template.dot.

Authors' Addresses

Joe Touch
USC/ISI
4676 Admiralty Way
Marina del Rey, CA 90292-6695
U.S.A.

Phone: +1 (310) 448-9151
Email: touch@isi.edu

W. Mark Townsley
Cisco
L'Atlantis, 11, Rue Camille Desmoulins
Issy Les Moulineaux, ILE DE FRANCE 92782

Email: townsley@cisco.com

APPENDIX A: Fragmentation efficiency

[A.1. Selecting fragment sizes](#)

There are different ways to fragment a packet. Consider a network with an MTU as shown in Figure 13, where packets are encapsulated over the same network layer as they arrive on (e.g., IP in IP). If a packet as large as the MTU arrives, it must be fragmented to accommodate the additional header.

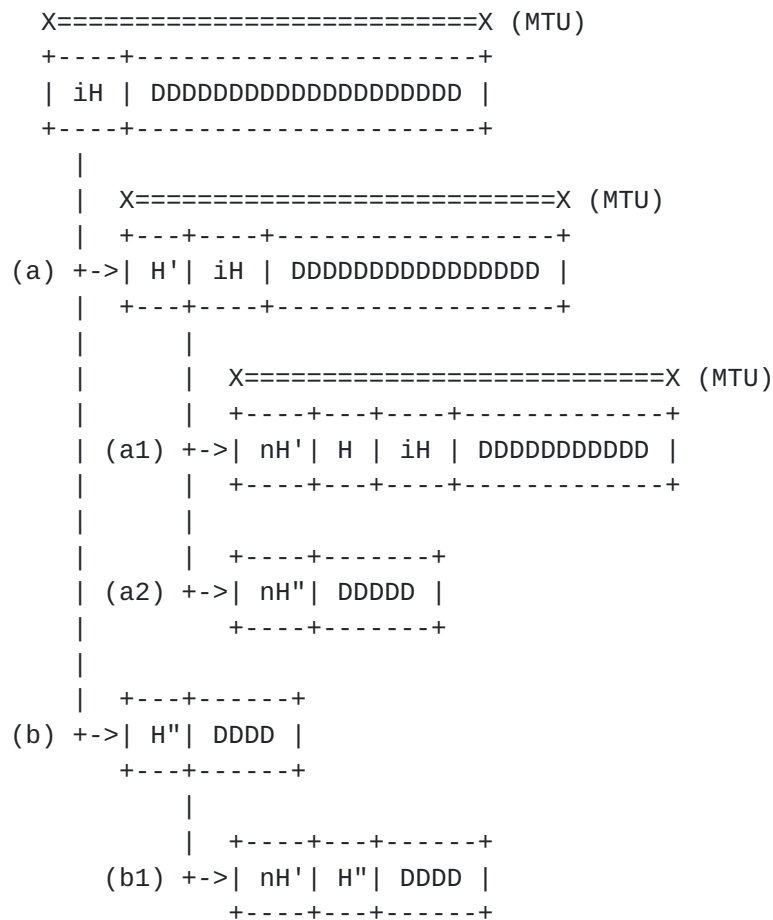


Figure 13 Fragmenting via maximum fit

Figure 13 shows this process, using Outer Fragmentation as an example (the situation is the same for Inner Fragmentation, but the headers that are affected differ). The arriving packet is first split into (a) and (b), where (a) is of the MTU of the network. However, this tunnel then traverses over another tunnel, whose impact the first tunnel ingress has not accommodated. The packet (a) arrives at the second tunnel ingress, and needs to be encapsulated again, but because it is already at the MTU, it needs to be fragmented as well,

into (a1) and (a2). In this case, packet (b) arrives at the second tunnel ingress and is encapsulated into (b1) without fragmentation, because it is already below the MTU size.

In Figure 14, the fragmentation is done evenly, i.e., by splitting the original packet into two roughly equal-sized components, (c) and (d). Note that (d) contains more packet data, because (c) includes the original packet header because this is an example of Outer Fragmentation. The packets (c) and (d) arrive at the second tunnel encapsulator, and are encapsulated again; this time, neither packet exceeds the MTU, and neither requires further fragmentation.

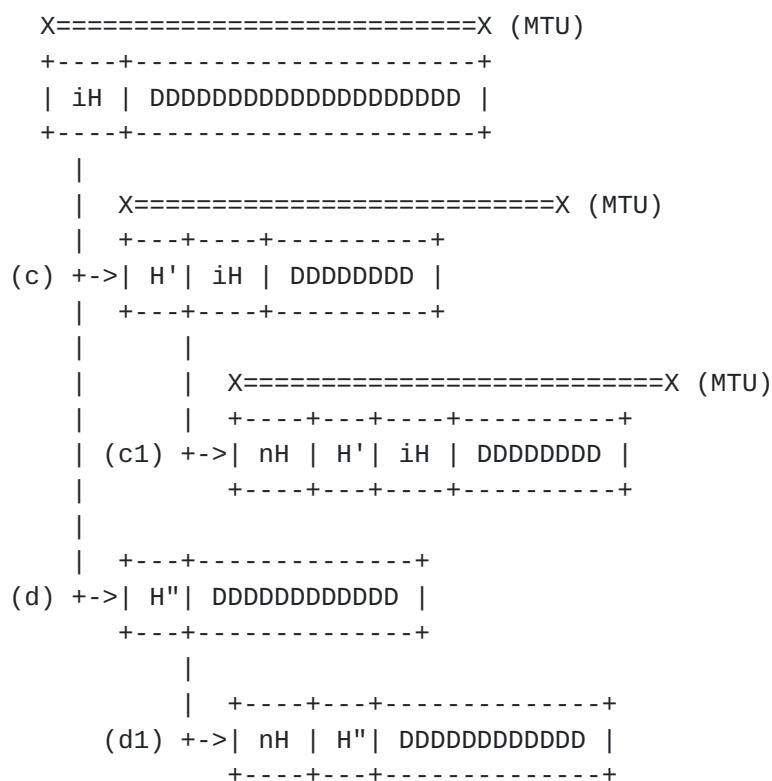


Figure 14 Fragmenting evenly

A.2. Packing

Encapsulating individual packets to traverse a tunnel can be inefficient, especially where headers are large relative to the packets being carried. In that case, it can be more efficient to encapsulate many small packets in a single, larger tunnel payload. This technique, similar to the effect of packet bursting in Gigabit Ethernet (regardless of whether they're encoded using L2 symbols as delineators), reduces the overhead of the encapsulation headers

(Figure 15). It reduces the work of header addition and removal at the tunnel endpoints, but increases other work involving the packing and unpacking of the component packets carried.

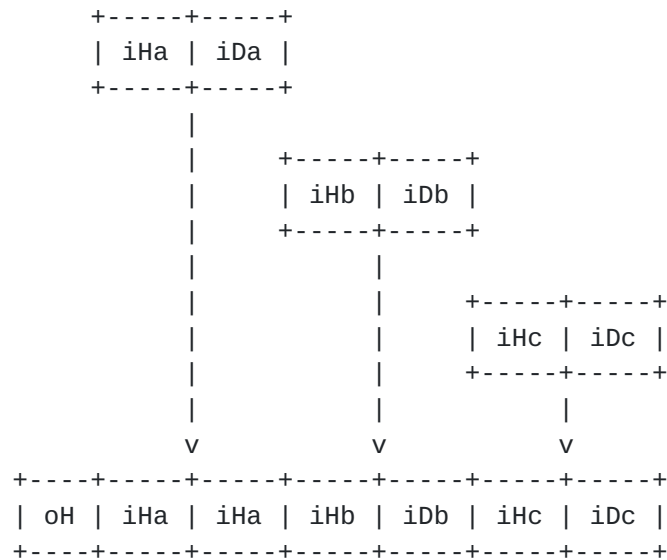


Figure 15 Packing packets into a tunnel

[NOTE: PPP chopping and coalescing?]