Authors: G. Fairhurst          T. Jones
         University of Aberdeen   University of Aberdeen
         M. Tuexen
         Muenster University of Applied Sciences
         I. Ruengeler
         Muenster University of Applied Sciences
         T. Voelker
         Muenster University of Applied Sciences

## Packetization Layer Path MTU Discovery for Datagram Transports

## Abstract

This document describes a robust method for Path MTU Discovery
(PMTUD) for datagram Packetization Layers (PLs). It describes an
extension to RFC 1191 and RFC 8201, which specifies ICMP-based Path
MTU Discovery for IPv4 and IPv6. The method allows a PL, or a
datagram application that uses a PL, to discover whether a network
path can support the current size of datagram. This can be used to
detect and reduce the message size when a sender encounters a
network black hole (where packets are discarded). The method can
probe a network path with progressively larger packets to discover
whether the maximum packet size can be increased. This allows a
sender to determine an appropriate packet size, providing
functionally for datagram transports that is equivalent to the
Packetization Layer PMTUD specification for TCP, specified in RFC
4821.

The document also provides implementation notes for incorporating
Datagram PMTUD into IETF datagram transports or applications that
use datagram transports.

When published, this specification updates RFC 4821.

## Status of This Memo

**Table of Contents**

## 1.  Introduction

The IETF has specified datagram transport using UDP, SCTP, and DCCP,
as well as protocols layered on top of these transports (e.g., SCTP/
UDP, DCCP/UDP, QUIC/UDP), and direct datagram transport over the IP
network layer. This document describes a robust method for Path MTU
Discovery (PMTUD) that may be used with these transport protocols
(or the applications that use their transport service) to discover
an appropriate size of packet to use across an Internet path.

### 1.1.  Classical Path MTU Discovery

Classical Path Maximum Transmission Unit Discovery (PMTUD) can be
used with any transport that is able to process ICMP Packet Too Big
(PTB) messages (e.g., [RFC1191] and [RFC8201]). In this document,
the term PTB message is applied to both IPv4 ICMP Unreachable
messages (type 3) that carry the error Fragmentation Needed (Type 3,
Code 4) [RFC0792] and ICMPv6 Packet Too Big messages (Type 2)
[RFC4443]. When a sender receives a PTB message, it reduces the
effective MTU to the value reported as the Link MTU in the PTB
message, and a method that from time-to-time increases the packet
size in attempt to discover an increase in the supported PMTU. The
packets sent with a size larger than the current effective PMTU are
known as probe packets.

Packets not intended as probe packets are either fragmented to the
current effective PMTU, or the attempt to send fails with an error
code. Applications are sometimes provided with a primitive to let
them read the Maximum Packet Size (MPS), derived from the current
effective PMTU.

Classical PMTUD is subject to protocol failures. One failure arises
when traffic using a packet size larger than the actual PMTU is
black-holed (all datagrams sent with this size, or larger, are
discarded). This could arise when the PTB messages are not delivered
back to the sender for some reason (see for example [RFC2923]).

Examples where PTB messages are not delivered include:

  *The generation of ICMP messages is usually rate limited. This
   could result in no PTB messages being generated to the sender
   (see section 2.4 of [RFC4443])

  *ICMP messages can be filtered by middleboxes (including
   firewalls) [RFC4890]. A stateful firewall could be configured
   with a policy to block incoming ICMP messages, which would
   prevent reception of PTB messages to a sending endpoint behind
   this firewall.

  *When the router issuing the ICMP message drops a tunneled packet,
   the resulting ICMP message will be directed to the tunnel
   ingress. This tunnel endpoint is responsible for forwarding the
   ICMP message and also processing the quoted packet within the
   payload field to remove the effect of the tunnel, and return a
   correctly formatted ICMP message to the sender [I-D.ietf-intarea-
   tunnels]. Failure to do this prevents the PTB message reaching
   the original sender.

  *Asymmetry in forwarding can result in there being no return route
   to the original sender, which would prevent an ICMP message being
   delivered to the sender. This issue can also arise when policy-
   based routing is used, Equal Cost Multipath (ECMP) routing is
   used, or a middlebox acts as an application load balancer. An
   example is where the path towards the server is chosen by ECMP
   routing depending on bytes in the IP payload. In this case, when
   a packet sent by the server encounters a problem after the ECMP
   router, then any resulting ICMP message needs to also be directed
   by the ECMP router towards the original sender.

  *There are additional cases where the next hop destination fails
   to receive a packet because of its size. This could be due to
   misconfiguration of the layer 2 path between nodes, for instance
   the MTU configured in a layer 2 switch, or misconfiguration of
   the Maximum Receive Unit (MRU). If the packet is dropped by the

link, this will not cause a PTB message to be sent to the
original sender.

Another failure could result if a node that is not on the network
path sends a PTB message that attempts to force a sender to change
the effective PMTU [RFC8201]. A sender can protect itself from
reacting to such messages by utilising the quoted packet within a
PTB message payload to validate that the received PTB message was
generated in response to a packet that had actually originated from
the sender. However, there are situations where a sender would be
unable to provide this validation. Examples where validation of the
PTB message is not possible include:

  *When a router issuing the ICMP message implements RFC792
   [RFC0792], it is only required to include the first 64 bits of
   the IP payload of the packet within the quoted payload. There
   could be insufficient bytes remaining for the sender to interpret
   the quoted transport information.

   Note: The recommendation in RFC1812 [RFC1812] is that IPv4
   routers return a quoted packet with as much of the original
   datagram as possible without the length of the ICMP datagram
   exceeding 576 bytes. IPv6 routers include as much of the invoking
   packet as possible without the ICMPv6 packet exceeding 1280 bytes
   [RFC4443].

  *The use of tunnels/encryption can reduce the size of the quoted
   packet returned to the original source address, increasing the
   risk that there could be insufficient bytes remaining for the
   sender to interpret the quoted transport information.

  *Even when the PTB message includes sufficient bytes of the quoted
   packet, the network layer could lack sufficient context to
   validate the message, because validation depends on information
   about the active transport flows at an endpoint node (e.g., the
   socket/address pairs being used, and other protocol header
   information).

  *When a packet is encapsulated/tunneled over an encrypted
   transport, the tunnel/encapsulation ingress might have
   insufficient context, or computational power, to reconstruct the
   transport header that would be needed to perform validation.

## 1.2.  Packetization Layer Path MTU Discovery

The term Packetization Layer (PL) has been introduced to describe
the layer that is responsible for placing data blocks into the
payload of IP packets and selecting an appropriate MPS. This
function is often performed by a transport protocol, but can also be

performed by other encapsulation methods working above the transport layer.

In contrast to PMTUD, Packetization Layer Path MTU Discovery (PLPMTUD) [RFC4821] does not rely upon reception and validation of PTB messages. It is therefore more robust than Classical PMTUD. This has become the recommended approach for implementing PMTU discovery.

It uses a general strategy where the PL sends probe packets to search for the largest size of unfragmented datagram that can be sent over a network path. Probe packets are sent with a progressively larger packet size. If a probe packet is successfully delivered (as determined by the PL), then the PLPMTU is raised to the size of the successful probe. If no response is received to a probe packet, the method reduces the probe size. The result of probing with the PLPMTU is used to set the application MPS.

PLPMTUD introduces flexibility in the implementation of PMTU discovery. At one extreme, it can be configured to only perform ICMP Black Hole Detection and recovery to increase the robustness of Classical PMTUD, or at the other extreme, all PTB processing can be disabled and PLPMTUD can completely replace Classical PMTUD (see Section 4.5).

PLPMTUD can also include additional consistency checks without increasing the risk that data is lost when probing to discover the path MTU. For example, information available at the PL, or higher layers, enables received PTB messages to be validated before being utilized.

1.3.  Path MTU Discovery for Datagram Services

Section 5 of this document presents a set of algorithms for datagram protocols to discover the largest size of unfragmented datagram that can be sent over a network path. The method described relies on features of the PL described in Section 3 and applies to transport protocols operating over IPv4 and IPv6. It does not require cooperation from the lower layers, although it can utilize PTB messages when these received messages are made available to the PL.

The UDP Usage Guidelines [RFC8085] state "an application SHOULD either use the Path MTU information provided by the IP layer or implement Path MTU Discovery (PMTUD)", but does not provide a mechanism for discovering the largest size of unfragmented datagram that can be used on a network path. Prior to this document, PLPMTUD had not been specified for UDP.

Section 10.2 of [RFC4821] recommends a PLPMTUD probing method for the Stream Control Transport Protocol (SCTP). SCTP utilizes probe packets consisting of a minimal sized HEARTBEAT chunk bundled with a

PAD chunk as defined in [RFC4820], but RFC4821 does not provide a complete specification. The present document provides the details to complete that specification.

The Datagram Congestion Control Protocol (DCCP) [RFC4340] requires implementations to support Classical PMTUD and states that a DCCP sender "MUST maintain the MPS allowed for each active DCCP session". It also defines the current congestion control MPS (CCMPS) supported by a network path. This recommends use of PMTUD, and suggests use of control packets (DCCP-Sync) as path probe packets, because they do not risk application data loss. The method defined in this specification could be used with DCCP.

Section 6 specifies the method for a set of transports, and provides information to enable the implementation of PLPMTUD with other datagram transports and applications that use datagram transports.

## 2.  Terminology

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

Other terminology is directly copied from [RFC4821], and the definitions in [RFC1122].

**Actual PMTU:**  The Actual PMTU is the PMTU of a network path between a sender PL and a destination PL, which the DPLPMTUD algorithm seeks to determine.

**Black Hole:**  A Black Hole is encountered when a sender is unaware that packets are not being delivered to the destination end point. Two types of Black Hole are relevant to DPLPMTUD:

   **Packet Black Hole:**  Packets encounter a Packet Black Hole when packets are not delivered to the destination endpoint (e.g., when the sender transmits packets of a particular size with a previously known effective PMTU and they are discarded by the network).

   **ICMP Black Hole**  An ICMP Black Hole is encountered when the sender is unaware that packets are not delivered to the destination endpoint because PTB messages are not received by the originating PL sender.

**Black holed :**  Traffic is black-holed when the sender is unaware that packets are not being delivered. This could be due to a Packet Black Hole or an ICMP Black Hole.

**Classical Path MTU Discovery:**
   Classical PMTUD is a process
   described in [RFC1191] and [RFC8201], in which nodes rely on PTB
   messages to learn the largest size of unfragmented datagram that
   can be used across a network path.

**Datagram:**  A datagram is a transport-layer protocol data unit,
   transmitted in the payload of an IP packet.

**Effective PMTU:**  The Effective PMTU is the current estimated value
   for PMTU that is used by a PMTUD. This is equivalent to the
   PLPMTU derived by PLPMTUD.

**EMTU_S:**  The Effective MTU for sending (EMTU_S) is defined in
   [RFC1122] as "the maximum IP datagram size that may be sent, for
   a particular combination of IP source and destination
   addresses...".

**EMTU_R:**  The Effective MTU for receiving (EMTU_R) is designated in
   [RFC1122] as the largest datagram size that can be reassembled by
   EMTU_R (Effective MTU to receive).

**Link:**  A Link is a communication facility or medium over which nodes
   can communicate at the link layer, i.e., a layer below the IP
   layer. Examples are Ethernet LANs and Internet (or higher) layer
   and tunnels.

**Link MTU:**  The Link Maximum Transmission Unit (MTU) is the size in
   bytes of the largest IP packet, including the IP header and
   payload, that can be transmitted over a link. Note that this
   could more properly be called the IP MTU, to be consistent with
   how other standards organizations use the acronym. This includes
   the IP header, but excludes link layer headers and other framing
   that is not part of IP or the IP payload. Other standards
   organizations generally define the link MTU to include the link
   layer headers.

**MAX_PMTU:**  The MAX_PMTU is the largest size of PLPMTU that DPLPMTUD
   will attempt to use.

**MPS:**  The Maximum Packet Size (MPS) is the largest size of
   application data block that can be sent across a network path by
   a PL. In DPLPMTUD this quantity is derived from the PLPMTU by
   taking into consideration the size of the lower protocol layer
   headers. Probe packets generated by DPLPMTUD can have a size
   larger than the MPS.

**MIN_PMTU:**  The MIN_PMTU is the smallest size of PLPMTU that DPLPMTUD
   will attempt to use.

**Packet:**
    A Packet is the IP header plus the IP payload.

**Packetization Layer (PL):**  The Packetization Layer (PL) is the layer
    of the network stack that places data into packets and performs
    transport protocol functions.

**Path:**  The Path is the set of links and routers traversed by a
    packet between a source node and a destination node by a
    particular flow.

**Path MTU (PMTU):**  The Path MTU (PMTU) is the minimum of the Link MTU
    of all the links forming a network path between a source node and
    a destination node.

**PTB_SIZE:**  The PTB_SIZE is a value reported in a validated PTB
    message that indicates next hop link MTU of a router along the
    path.

**PLPMTU:**  The Packetization Layer PMTU is an estimate of the actual
    PMTU provided by the DPLPMTUD algorithm.

**PLPMTUD:**  Packetization Layer Path MTU Discovery (PLPMTUD), the
    method described in this document for datagram PLs, which is an
    extension to Classical PMTU Discovery.

**Probe packet:**  A probe packet is a datagram sent with a purposely
    chosen size (typically the current PLPMTU or larger) to detect if
    packets of this size can be successfully sent end-to-end across
    the network path.

3.  **Features Required to Provide Datagram PLPMTUD**

   TCP PLPMTUD has been defined using standard TCP protocol mechanisms.
   All of the requirements in [RFC4821] also apply to the use of the
   technique with a datagram PL. Unlike TCP, some datagram PLs require
   additional mechanisms to implement PLPMTUD.

   There are eight requirements for performing the datagram PLPMTUD
   method described in this specification:

   1. PMTU parameters: A DPLPMTUD sender is RECOMMENDED to provide
      information about the maximum size of packet that can be
      transmitted by the sender on the local link (the local Link
      MTU). It MAY utilize similar information about the receiver
      when this is supplied (note this could be less than EMTU_R).
      This avoids implementations trying to send probe packets that
      can not be transmitted by the local link. Too high of a value
      could reduce the efficiency of the search algorithm. Some
      applications also have a maximum transport protocol data unit

(PDU) size, in which case there is no benefit from probing for a size larger than this (unless a transport allows multiplexing multiple applications PDUs into the same datagram).

2. PLPMTU: A datagram application using a PL not supporting fragmentation is REQUIRED to be able to choose the size of datagrams sent to the network, up to the PLPMTU, or a smaller value (such as the MPS) derived from this. This value is managed by the DPLPMTUD method. The PLPMTU (specified as the effective PMTU in Section 1 of [RFC1191]) is equivalent to the EMTU_S (specified in [RFC1122]).

3. Probe packets: On request, a DPLPMTUD sender is REQUIRED to be able to transmit a packet larger than the PLMPMTU. This is used to send a probe packet. In IPv4, a probe packet MUST be sent with the Don't Fragment (DF) bit set in the IP header, and without network layer endpoint fragmentation. In IPv6, a probe packet is always sent without source fragmentation (as specified in section 5.4 of [RFC8201]).

4. Processing PTB messages: A DPLPMTUD sender MAY optionally utilize PTB messages received from the network layer to help identify when a network path does not support the current size of probe packet. Any received PTB message MUST be validated before it is used to update the PLPMTU discovery information [RFC8201]. This validation confirms that the PTB message was sent in response to a packet originating by the sender, and needs to be performed before the PLPMTU discovery method reacts to the PTB message. A PTB message MUST NOT be used to increase the PLPMTU [RFC8201].

5. Reception feedback: The destination PL endpoint is REQUIRED to provide a feedback method that indicates to the DPLPMTUD sender when a probe packet has been received by the destination PL endpoint. The mechanism needs to be robust to the possibility that packets could be significantly delayed along a network path. The local PL endpoint at the sending node is REQUIRED to pass this feedback to the sender DPLPMTUD method.

6. Probe loss recovery: It is RECOMMENDED to use probe packets that do not carry any user data that would require retransmission if lost. Most datagram transports permit this. If a probe packet contains user data requiring retransmission in case of loss, the PL (or layers above) are REQUIRED to arrange any retransmission/repair of any resulting loss. DPLPMTUD is REQUIRED to be robust in the case where probe packets are lost due to other reasons (including link transmission error, congestion).

7. Probing and congestion control: The DPLPMTUD sender treats isolated loss of a probe packet (with or without a corresponding PTB message) as a potential indication of a PMTU limit for the path. Loss of a probe packet SHOULD NOT be treated as an indication of congestion. The loss of a probe packet SHOULD NOT directly trigger a congestion control reaction [RFC4821] because this could result in unecessary reduction of the sending rate. The interval between probe packets MUST be at least one RTT.

8. Shared PLPMTU state: The PLPMTU value MAY also be stored with the corresponding entry associated with the destination in the IP layer cache, and used by other PL instances. The specification of PLPMTUD [RFC4821] states: "If PLPMTUD updates the MTU for a particular path, all Packetization Layer sessions that share the path representation (as described in Section 5.2 of [RFC4821]) SHOULD be notified to make use of the new MTU". Such methods MUST be robust to the wide variety of underlying network forwarding behaviors. Section 5.2 of [RFC8201] provides guidance on the caching of PMTU information and also the relation to IPv6 flow labels.

In addition, the following principles are stated for design of a DPLPMTUD method:

  *MPS: A method is REQUIRED to signal an appropriate MPS to the higher layer using the PL. The value of the MPS can change following a change to the path. It is RECOMMENDED that methods avoid forcing an application to use an arbitrary small MPS (PLPMTU) for transmission while the method is searching for the currently supported PLPMTU. Datagram PLs do not necessarily support fragmentation of PDUs larger than the PLPMTU. A reduced MPS can adversely impact the performance of a datagram application.

  *Path validation: It is RECOMMENDED that methods are robust to path changes that could have occurred since the path characteristics were last confirmed, and to the possibility of inconsistent path information being received.

  *Datagram reordering: A method is REQUIRED to be robust to the possibility that a flow encounters reordering, or the traffic (including probe packets) is divided over more than one network path.

  *When to probe: It is RECOMMENDED that methods determine whether the path has changed since it last measured the path. This can help determine when to probe the path again.

4.  DPLPMTUD Mechanisms

   This section lists the protocol mechanisms used in this
   specification.

4.1.  PLPMTU Probe Packets

   The DPLPMTUD method relies upon the PL sender being able to generate
   probe packets with a specific size. TCP is able to generate these
   probe packets by choosing to appropriately segment data being sent
   [RFC4821]. In contrast, a datagram PL that needs to construct a
   probe packet has to either request an application to send a data
   block that is larger than that generated by an application, or to
   utilize padding functions to extend a datagram beyond the size of
   the application data block. Protocols that permit exchange of
   control messages (without an application data block) MAY prefer to
   generate a probe packet by extending a control message with padding
   data.

   A receiver is REQUIRED to be able to distinguish an in-band data
   block from any added padding. This is needed to ensure that any
   added padding is not passed on to an application at the receiver.

   This results in three possible ways that a sender can create a probe
   packet:

   Probing using padding data:  A probe packet that contains only
      control information together with any padding, which is needed to
      be inflated to the size required for the probe packet. Since
      these probe packets do not carry an application-supplied data
      block, they do not typically require retransmission, although
      they do still consume network capacity and incur endpoint
      processing.

   Probing using application data and padding data:  A probe packet
      that contains a data block supplied by an application that is
      combined with padding to inflate the length of the datagram to
      the size required for the probe packet. If the application/
      transport needs protection from the loss of this probe packet,
      the application/transport could perform transport-layer
      retransmission/repair of the data block (e.g., by retransmission
      after loss is detected or by duplicating the data block in a
      datagram without the padding data).

   Probing using application data:  A probe packet that contains a data
      block supplied by an application that matches the size required
      for the probe packet. This method requests the application to
      issue a data block of the desired probe size. If the application/
      transport needs protection from the loss of an unsuccessful probe
      packet, the application/transport needs then to perform

transport-layer retransmission/repair of the data block (e.g., by
retransmission after loss is detected).

A PL that uses a probe packet carrying an application data block,
could need to retransmit this application data block if the probe
fails. This could need the PL to re-fragment the data block to a
smaller packet size that is expected to traverse the end-to-end path
(which could utilize endpoint network-layer or PL fragmentation when
these are available).

DPLPMTUD MAY choose to use only one of these methods to simplify the
implementation.

Probe messages sent by a PL MUST contain enough information to
uniquely identify the probe within Maximum Segment Lifetime, while
being robust to reordering and replay of probe response and PTB
messages.

## 4.2.  Confirmation of Probed Packet Size

The PL needs a method to determine (confirm) when probe packets have
been successfully received end-to-end across a network path.

Transport protocols can include end-to-end methods that detect and
report reception of specific datagrams that they send (e.g., DCCP
and SCTP provide keep-alive/heartbeat features). When supported,
this mechanism SHOULD also be used by DPLPMTUD to acknowledge
reception of a probe packet.

A PL that does not acknowledge data reception (e.g., UDP and UDP-
Lite) is unable itself to detect when the packets that it sends are
discarded because their size is greater than the actual PMTU. These
PLs need to either rely on an application protocol to detect this
loss.

Section 6 specifies this function for a set of IETF-specified
protocols.

## 4.3.  Detection of Unsupported PLPMTU Size, aka Black Hole Detection

A PL sender needs to reduce the PLPMTU when it discovers the actual
PMTU supported by a network path is less than the PLPMTU. This can
be triggered when a validated PTB message is received, or by another
event that indicates the network path no longer sustains the current
packet size, such as a loss report from the PL, or repeated lack of
response to probe packets sent to confirm the PLPMTU. Detection is
followed by a reduction of the PLPMTU.

This is performed by sending packet probes of size PLPMTU to verify
that a network path still supports the last acknowledged PLPMTU
size. There are two alternative mechanism:

  *A PL can rely upon a mechanism implemented within the PL to
   detect excessive loss of data sent with a specific packet size
   and then conclude that this excessive loss could be a result of
   an invalid PMTU (as in PLPMTUD for TCP [RFC4821]).

  *A PL can use the DPLPMTUD probing mechanism to periodically
   generate probe packets of the size of the current PLPMTU (e.g.,
   using the confirmation timer Section 5.1.1). A timer tracks
   whether acknowledgments are received. Successive loss of probes
   is an indication that the current path no longer supports the
   PLPMTU (e.g., when the number of probe packets sent without
   receiving an acknowledgement, PROBE_COUNT, becomes greater than
   MAX_PROBES).

A PL MAY inhibit sending probe packets when no application data has
been sent since the previous probe packet. A PL preferring to use an
up-to-data PLPMTU once user data is sent again, MAY choose to
continue PLPMTU discovery for each path. However, this may result in
additional packets being sent.

When the method detects the current PLPMTU is not supported,
DPLPMTUD sets a lower MPS. The PL then confirms that the updated
PLPMTU can be successfully used across the path. The PL could need
to send a probe packet with a size less than the size of the data
block generated by an application. In this case, the PL could
provide a way to fragment a datagram at the PL, or use a control
packet as the packet probe.

## 4.4.  Disabling the Effect of PMTUD

A PL implementing this specification MUST suspend network layer
processing of outgoing packets that enforces a PMTU [RFC1191]
[RFC8201] for each flow utilising DPLPMTUD, and instead use DPLPMTUD
to control the size of packets that are sent by a flow. This removes
the need for the network layer to drop or fragment sent packets that
have a size greater than the PMTU.

## 4.5.  Response to PTB Messages

This method requires the DPLPMTUD sender to validate any received
PTB message before using the PTB information. The response to a PTB
message depends on the PTB_SIZE indicated in the PTB message, the
state of the PLPMTUD state machine, and the IP protocol being used.

Section 4.5.1 first describes validation for both IPv4 ICMP
Unreachable messages (type 3) and ICMPv6 Packet Too Big messages,
both of which are referred to as PTB messages in this document.

## 4.5.1.  Validation of PTB Messages

This section specifies utilization of PTB messages.

  *A simple implementation MAY ignore received PTB messages and in
   this case the PLPMTU is not updated when a PTB message is
   received.

  *An implementation that supports PTB messages MUST validate
   messages before they are further processed.

A PL that receives a PTB message from a router or middlebox,
performs ICMP validation as specified in Section 5.2 of [RFC8085]
[RFC8201]. Because DPLPMTUD operates at the PL, the PL needs to
check that each received PTB message is received in response to a
packet transmitted by the endpoint PL performing DPLPMTUD.

The PL MUST check the protocol information in the quoted packet
carried in an ICMP PTB message payload to validate the message
originated from the sending node. This validation includes
determining that the combination of the IP addresses, the protocol,
the source port and destination port match those returned in the
quoted packet - this is also necessary for the PTB message to be
passed to the corresponding PL.

The validation SHOULD utilize information that it is not simple for
an off-path attacker to determine [RFC8085]. For example, by
checking the value of a protocol header field known only to the two
PL endpoints. A datagram application that uses well-known source and
destination ports ought to also rely on other information to
complete this validation.

These checks are intended to provide protection from packets that
originate from a node that is not on the network path. A PTB message
that does not complete the validation MUST NOT be further utilized
by the DPLPMTUD method.

PTB messages that have been validated MAY be utilized by the
DPLPMTUD algorithm, but MUST NOT be used directly to set the PLPMTU.
A method that utilizes these PTB messages can improve the speed at
the which the algorithm detects an appropriate PLPMTU, compared to
one that relies solely on probing. Section 4.5.2 describes this
processing.

### 4.5.2.  Use of PTB Messages

A set of checks are intended to provide protection from a router
that reports an unexpected PTB_SIZE. The PL also needs to check that
the indicated PTB_SIZE is less than the size used by probe packets
and larger than minimum size accepted.

This section provides a summary of how PTB messages can be utilized.
This processing depends on the PTB_SIZE and the current value of a
set of variables:

**PTB_SIZE < MIN_PMTU**
>       *Invalid PTB_SIZE see [Section 4.5.1](#).
>
>       *PTB message ought to be discarded without further
>        processing (e. g. PLPMTU not modified).
>
>       *The information could be utilized as an input to trigger
>        enabling a resilience mode.

**MIN_PMTU < PTB_SIZE < BASE_PMTU**
>       *A robust PL MAY enter an error state (see [Section 5.2](#)) for
>        an IPv4 path when the PTB_SIZE reported in the PTB message
>        is larger than or equal to 68 bytes [[RFC0791](#)] and when this
>        is less than the BASE_PMTU.
>
>       *A robust PL MAY enter an error state (see [Section 5.2](#)) for
>        an IPv6 path when the PTB_SIZE reported in the PTB message
>        is larger than or equal to 1280 bytes [[RFC8200](#)] and when
>        this is less than the BASE_PMTU.

**PTB_SIZE = PLPMTU**
>       *Completes the search for a larger PLPMTU.

**PTB_SIZE > PROBED_SIZE**
>       *Inconsistent network signal.
>
>       *PTB message ought to be discarded without further
>        processing (e. g. PLPMTU not modified).
>
>       *The information could be utilized as an input to trigger
>        enabling a resilience mode.

**BASE_PMTU <= PTB_SIZE < PLPMTU**
>       *Black Hole Detection is triggered and the PLPMTU ought to
>        be set to BASE_PMTU.
>
>       *The PL could use the PTB_SIZE reported in the PTB message
>        to initialize a search algorithm.

**PLPMTU < PTB_SIZE < PROBED_SIZE**

> *The PLPMTU continues to be valid, but the last PROBED_SIZE searched was larger than the actual PMTU.

> *The PLPMTU is not updated.

> *The PL can use the reported PTB_SIZE from the PTB message as the next search point when it resumes the search algorithm.

## 5. Datagram Packetization Layer PMTUD

This section specifies Datagram PLPMTUD (DPLPMTUD). The method can be introduced at various points (as indicated with * in the figure below) in the IP protocol stack to discover the PLPMTU so that an application can utilize an appropriate MPS for the current network path. DPLPMTUD SHOULD NOT be used by an application if it is already used in a lower layer.



Figure 1: Examples where DPLPMTUD can be implemented

The central idea of DPLPMTUD is probing by a sender. Probe packets are sent to find the maximum size of a user message that can be completely transferred across the network path from the sender to the destination.

The following sections identify the components needed for implementation, provides an overview of the phases of operation, and specifies the state machine and search algorithm.

### 5.1.  DPLPMTUD Components

This section describes the timers, constants, and variables of
DPLPMTUD.

### 5.1.1.  Timers

The method utilizes up to three timers:

**PROBE_TIMER:**  The PROBE_TIMER is configured to expire after a period
   longer than the maximum time to receive an acknowledgment to a
   probe packet. This value MUST NOT be smaller than 1 second, and
   SHOULD be larger than 15 seconds. Guidance on selection of the
   timer value are provided in section 3.1.1 of the UDP Usage
   Guidelines [RFC8085].

   If the PL has a path Round Trip Time (RTT) estimate and timely
   acknowledgements the PROBE_TIMER can be derived from the PL RTT
   estimate.

**PMTU_RAISE_TIMER:**  The PMTU_RAISE_TIMER is configured to the period
   a sender will continue to use the current PLPMTU, after which it
   re-enters the Search phase. This timer has a period of 600
   seconds, as recommended by PLPMTUD [RFC4821].

   DPLPMTUD MAY inhibit sending probe packets when no application
   data has been sent since the previous probe packet. A PL
   preferring to use an up-to-data PMTU once user data is sent
   again, can choose to continue PMTU discovery for each path.
   However, this may result in sending additional packets.

**CONFIRMATION_TIMER:**  When an acknowledged PL is used, this timer
   MUST NOT be used. For other PLs, the CONFIRMATION_TIMER is
   configured to the period a PL sender waits before confirming the
   current PLPMTU is still supported. This is less than the
   PMTU_RAISE_TIMER and used to decrease the PLPMTU (e.g., when a
   black hole is encountered). Confirmation needs to be frequent
   enough when data is flowing that the sending PL does not black
   hole extensive amounts of traffic. Guidance on selection of the
   timer value are provided in section 3.1.1 of the UDP Usage
   Guidelines [RFC8085].

   DPLPMTUD MAY inhibit sending probe packets when no application
   data has been sent since the previous probe packet. A PL
   preferring to use an up-to-data PMTU once user data is sent
   again, can choose to continue PMTU discovery for each path.
   However, this may result in sending additional packets.

An implementation could implement the various timers using a single
timer.

### 5.1.2.  Constants

The following constants are defined:

**MAX_PROBES:**  The MAX_PROBES is the maximum value of the PROBE_COUNT
   counter (see Section 5.1.3). MAX_PROBES represents the limit for
   the number of consecutive probe attempts of any size. The default
   value of MAX_PROBES is 3. This value is greater than 1 to provide
   robustness to isolated packet loss.

**MIN_PMTU:**  The MIN_PMTU is the smallest allowed probe packet size.
   For IPv6, this value is 1280 bytes, as specified in [RFC8200].
   For IPv4, the minimum value is 68 bytes.

   Note: An IPv4 router is required to be able to forward a datagram
   of 68 bytes without further fragmentation. This is the combined
   size of an IPv4 header and the minimum fragment size of 8 bytes.
   In addition, receivers are required to be able to reassemble
   fragmented datagrams at least up to 576 bytes, as stated in
   section 3.3.3 of [RFC1122].

**MAX_PMTU:**  The MAX_PMTU is the largest size of PLPMTU. This has to
   be less than or equal to the minimum of the local MTU of the
   outgoing interface and the destination PMTU for receiving. An
   application, or PL, MAY choose a smaller MAX_PMTU when there is
   no need to send packets larger than a specific size.

**BASE_PMTU:**  The BASE_PMTU is a configured size expected to work for
   most paths. The size is equal to or larger than the MIN_PMTU and
   smaller than the MAX_PMTU. In the case of IPv6, this value is
   1280 bytes [RFC8200]. When using IPv4, a size of 1200 bytes is
   RECOMMENDED.

### 5.1.3.  Variables

This method utilizes a set of variables:

**PROBED_SIZE:**  The PROBED_SIZE is the size of the current probe
   packet. This is a tentative value for the PLPMTU, which is
   awaiting confirmation by an acknowledgment.

**PROBE_COUNT:**  The PROBE_COUNT is a count of the number of successive
   unsuccessful probe packets that have been sent. Each time a probe
   packet is acknowledged, the value is set to zero.

The figure below illustrates the relationship between the packet
size constants and variables at a point of time when the DPLPMTUD
algorithm performs path probing to increase the size of the PLPMTU.
A probe packet has been sent of size PROBED_SIZE. Once this is
acknowledged, the PLPMTU will raise to PROBED_SIZE allowing the

DPLPMTUD algorithm to further increase PROBED_SIZE towards the
actual PMTU.

```
MIN PMTU                                          MAX PMTU
   <-------------------------------------------------->
                    |           |        |
                    v           |        v
              BASE PMTU         |    Actual PMTU
                                v
                           PROBED SIZE
                    |
                    v
                 PLPMTU
```

Figure 2: Relationships between packet size constants and variables

**5.1.4.  Overview of DPLPMTUD Phases**

This section provides a high-level informative view of the DPLPMTUD
method, by describing the movement of the method through several
phases of operation. More detail is available in the state machine
Section 5.2.

```
                  +------+      Connectivity
              +-->| Base |----- or BASE PMTU
              |   +------+      confirmation failed
              |      |                |
              |      |                v
              | Connectivity      +-------+
              | and BASE PMTU     | Error |
              | confirmed         +-------+
              |      |                |
              |      v           Consistent
  PLPMTU      | +--------+       connectivity
confirmation  +-| Search |<----- and BASE PMTU
  failed        +--------+       confirmed
              |   ^     |
              | Raise   Search
              | timer   algorithm
              | expired completed
              |   |     |
              |   |     v
              | +-----------------+
              +-| Search Complete |
                +-----------------+
```

Figure 3: DPLPMTUD Phases

**Base:**  The Base Phase confirms connectivity to the remote peer. This
phase is implicit for a connection-oriented PL (where it can be
performed in a PL connection handshake). A connectionless PL
needs to send an acknowledged probe packet to confirm that the

remote peer is reachable. The sender also confirms that BASE_PMTU
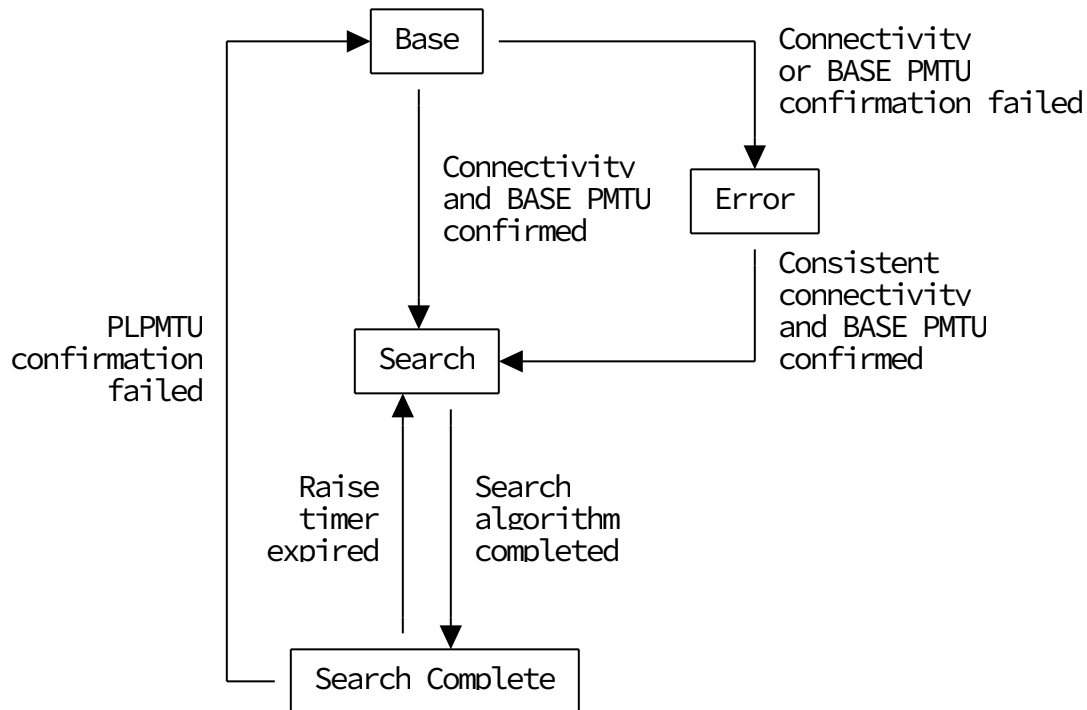is supported across the network path.

A PL that does not wish to support a path with a PLPMTU less than
BASE_PMTU can simplify the phase into a single step by performing
the connectivity checks with a probe of the BASE_PMTU size.

Once confirmed, DPLPMTUD enters the Search Phase. If this phase
fails to confirm, DPLPMTUD enters the Error Phase.

Search:  The Search Phase utilizes a search algorithm to send probe
packets to seek to increase the PLPMTU. The algorithm concludes
when it has found a suitable PLPMTU, by entering the Search
Complete Phase.

A PL could respond to PTB messages using the PTB to advance or
terminate the search, see Section 4.5.

Search Complete:  The Search Complete Phase is entered when the
PLPMTU is supported across the network path. A PL can use a
CONFIRMATION_TIMER to periodically repeat a probe packet for the
current PLPMTU size. If the sender is unable to confirm
reachability (e.g., if the CONFIRMATION_TIMER expires) or the PL
signals a lack of reachability, DPLPMTUD enters the Base phase.

The PMTU_RAISE_TIMER is used to periodically resume the search
phase to discover if the PLPMTU can be raised. Black Hole
Detection or receipt of a validated PTB message (see Section
4.5.1) can cause the sender to enter the Base Phase.

Error:  The Error Phase is entered when there is conflicting or
invalid PLPMTU information for the path (e.g. a failure to
support the BASE_PMTU) that cause DPLPMTUD to be unable to
progress and the PLPMTU is lowered.

DPLPMTUD remains in the Error Phase until a consistent view of
the path can be discovered and it has also been confirmed that
the path supports the BASE_PMTU (or DPLPMTUD is suspended).

An implementation that only reduces the PLPMTU to a suitable size
would be sufficient to ensure reliable operation, but can be very
inefficient when the actual PMTU changes or when the method (for
whatever reason) makes a suboptimal choice for the PLPMTU.

A full implementation of DPLPMTUD provides an algorithm enabling the
DPLPMTUD sender to increase the PLPMTU following a change in the
characteristics of the path, such as when a link is reconfigured
with a larger MTU, or when there is a change in the set of links
traversed by an end-to-end flow (e.g., after a routing or path fail-
over decision).

## 5.2.  State Machine

A state machine for DPLPMTUD is depicted in Figure 4. If multipath or multihoming is supported, a state machine is needed for each path.

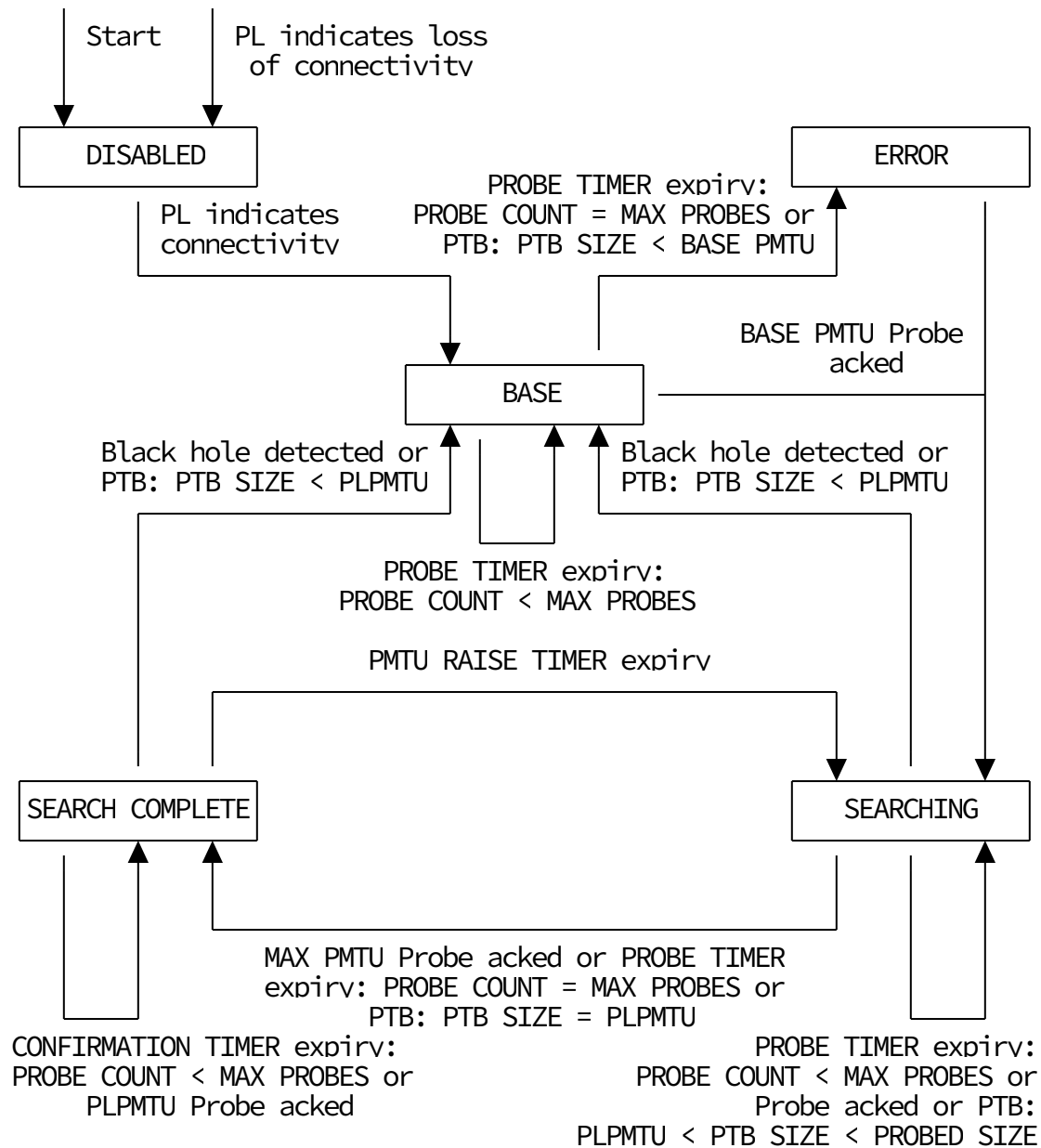Note: Not all changes are not shown to simplify the diagram.



Figure 4: State machine for Datagram PLPMTUD

The following states are defined:

**DISABLED:**
>        The DISABLED state is the initial state before probing
>    has started. It is also entered from any other state, when the PL
>    indicates loss of connectivity. This state is left, once the PL
>    indicates connectivity to the remote PL.

**BASE:**   The BASE state is used to confirm that the BASE_PMTU size is
>    supported by the network path and is designed to allow an
>    application to continue working when there are transient
>    reductions in the actual PMTU. It also seeks to avoid long
>    periods where traffic is black holed while searching for a larger
>    PLPMTU.

>    On entry, the PROBED_SIZE is set to the BASE_PMTU size and the
>    PROBE_COUNT is set to zero.

>    Each time a probe packet is sent, the PROBE_TIMER is started. The
>    state is exited when the probe packet is acknowledged, and the PL
>    sender enters the SEARCHING state.

>    The state is also left when the PROBE_COUNT reaches MAX_PROBES or
>    a received PTB message is validated. This causes the PL sender to
>    enter the ERROR state.

**SEARCHING:**   The SEARCHING state is the main probing state. This
>    state is entered when probing for the BASE_PMTU was successful.

>    Each time a probe packet is acknowledged, the PROBE_COUNT is set
>    to zero, the PLPMTU is set to the PROBED_SIZE and then the
>    PROBED_SIZE is increased using the search algorithm.

>    When a probe packet is sent and not acknowledged within the
>    period of the PROBE_TIMER, the PROBE_COUNT is incremented and a
>    new probe packet is transmitted. The state is exited when the
>    PROBE_COUNT reaches MAX_PROBES, a received PTB message is
>    validated, a probe of size MAX_PMTU is acknowledged, or a black
>    hole is detected.

**SEARCH_COMPLETE:**   The SEARCH_COMPLETE state indicates a successful
>    end to the SEARCHING state. DPLPMTUD remains in this state until
>    either the PMTU_RAISE_TIMER expires, a received PTB message is
>    validated, or a black hole is detected.

>    When DPLPMTUD uses an unacknowledged PL and is in the
>    SEARCH_COMPLETE state, a CONFIRMATION_TIMER periodically resets
>    the PROBE_COUNT and schedules a probe packet with the size of the
>    PLPMTU. If MAX_PROBES successive PLPMTUD sized probes fail to be
>    acknowledged the method enters the BASE state. When used with an
>    acknowledged PL (e.g., SCTP), DPLPMTUD SHOULD NOT continue to
>    generate PLPMTU probes in this state.

ERROR:

   The ERROR state represents the case where either the network
   path is not known to support a PLPMTU of at least the BASE_PMTU
   size or when there is contradictory information about the network
   path that would otherwise result in excessive variation in the
   MPS signalled to the higher layer. The state implements a method
   to mitigate oscillation in the state-event engine. It signals a
   conservative value of the MPS to the higher layer by the PL. The
   state is exited when packet probes no longer detect the error or
   when the PL indicates that connectivity has been lost.

   Implementations are permitted to enable endpoint fragmentation if
   the DPLPMTUD is unable to validate MIN_PMTU within PROBE_COUNT
   probes. If DPLPMTUD is unable to validate MIN_PMTU the
   implementation should transition to the DISABLED state.

   Note: MIN_PMTU may be identical to BASE_PMTU, simplifying the
   actions in this state.

## 5.3.  Search to Increase the PLPMTU

This section describes the algorithms used by DPLPMTUD to search for
a larger PLPMTU.

### 5.3.1.  Probing for a larger PLPMTU

Implementations use a search algorithm across the search range to
determine whether a larger PLPMTU can be supported across a network
path.

The method discovers the search range by confirming the minimum
PLPMTU and then using the probe method to select a PROBED_SIZE less
than or equal to MAX_PMTU. MAX_PMTU is the minimum of the local MTU
and EMTU_R (learned from the remote endpoint). The MAX_PMTU MAY be
reduced by an application that sets a maximum to the size of
datagrams it will send.

The PROBE_COUNT is initialized to zero when the first probe with a
size greater than or equal to PLPMTUD is sent. A timer is used by
the search algorithm to trigger the sending of probe packets of size
PROBED_SIZE, larger than the PLPMTU. Each probe packet successfully
sent to the remote peer is confirmed by acknowledgement at the PL,
see Section 4.1.

Each time a probe packet is sent to the destination, the PROBE_TIMER
is started. The timer is canceled when the PL receives
acknowledgment that the probe packet has been successfully sent
across the path Section 4.1. This confirms that the PROBED_SIZE is
supported, and the PROBED_SIZE value is then assigned to the PLPMTU.

The search algorithm can continue to send subsequent probe packets
of an increasing size.

If the timer expires before a probe packet is acknowledged, the
probe has failed to confirm the PROBED_SIZE. Each time the
PROBE_TIMER expires, the PROBE_COUNT is incremented, the PROBE_TIMER
is reinitialized, and a new probe of the same size or any other size
(determined by the search algorithm) can be sent. The maximum number
of consecutive failed probes is configured (MAX_PROBES). If the
value of the PROBE_COUNT reaches MAX_PROBES, probing will stop, and
the PL sender enters the SEARCH_COMPLETE state.

### 5.3.2.  Selection of Probe Sizes

The search algorithm needs to determine a minimum useful gain in
PLPMTU. It would not be constructive for a PL sender to attempt to
probe for all sizes. This would incur unnecessary load on the path
and has the undesirable effect of slowing the time to reach a more
optimal MPS. Implementations SHOULD select the set of probe packet
sizes to maximize the gain in PLPMTU from each search step.

Implementations could optimize the search procedure by selecting
step sizes from a table of common PMTU sizes. When selecting the
appropriate next size to search, an implementer ought to also
consider that there can be common sizes of MPS that applications
seek to use, and their could be common sizes of MTU used within the
network.

### 5.3.3.  Resilience to Inconsistent Path Information

A decision to increase the PLPMTU needs to be resilient to the
possibility that information learned about the network path is
inconsistent. A path is inconsistent, when, for example, probe
packets are lost due to other reasons (i.e. not packet size) or due
to frequent path changes. Frequent path changes could occur by
unexpected "flapping" - where some packets from a flow pass along
one path, but other packets follow a different path with different
properties.

A PL sender is able to detect inconsistency from the sequence of
PLPMTU probes that it sends or the sequence of PTB messages that it
receives. When inconsistent path information is detected, a PL
sender could use an alternate search mode that clamps the offered
MPS to a smaller value for a period of time. This avoids unnecessary
loss of packets due to MTU limitation.

### 5.4.  Robustness to Inconsistent Paths

Some paths could be unable to sustain packets of the BASE_PMTU size.
To be robust to these paths an implementation could implement the

Error State. This allows fallback to a smaller than desired PLPMTU, rather than suffer connectivity failure. This could utilize methods such as endpoint IP fragmentation to enable the PL sender to communicate using packets smaller than the BASE_PMTU.

## 6.  Specification of Protocol-Specific Methods

DPLPMTUD requires protocol-specific details to be specified for each PL that is used.

The first subsection provides guidance on how to implement the DPLPMTUD method as a part of an application using UDP or UDP-Lite. The guidance also applies to other datagram services that do not include a specific transport protocol (such as a tunnel encapsulation). The following subsections describe how DPLPMTUD can be implemented as a part of the transport service, allowing applications using the service to benefit from discovery of the PLPMTU without themselves needing to implement this method.

### 6.1.  Application support for DPLPMTUD with UDP or UDP-Lite

The current specifications of UDP [RFC0768] and UDP-Lite [RFC3828] do not define a method in the RFC-series that supports PLPMTUD. In particular, the UDP transport does not provide the transport layer features needed to implement datagram PLPMTUD.

The DPLPMTUD method can be implemented as a part of an application built directly or indirectly on UDP or UDP-Lite, but relies on higher-layer protocol features to implement the method [RFC8085].

Some primitives used by DPLPMTUD might not be available via the Datagram API (e.g., the ability to access the PLPMTU from the IP layer cache, or interpret received PTB messages).

In addition, it is desirable that PMTU discovery is not performed by multiple protocol layers. An application SHOULD avoid using DPLPMTUD when the underlying transport system provides this capability. To use common method for managing the PLPMTU has benefits, both in the ability to share state between different processes and opportunities to coordinate probing.

### 6.1.1.  Application Request

An application needs an application-layer protocol mechanism (such as a message acknowledgement method) that solicits a response from a destination endpoint. The method SHOULD allow the sender to check the value returned in the response to provide additional protection from off-path insertion of data [RFC8085], suitable methods include a parameter known only to the two endpoints, such as a session ID or initialized sequence number.

### 6.1.2.  Application Response

An application needs an application-layer protocol mechanism to
communicate the response from the destination endpoint. This
response may indicate successful reception of the probe across the
path, but could also indicate that some (or all packets) have failed
to reach the destination.

### 6.1.3.  Sending Application Probe Packets

A probe packet that may carry an application data block, but the
successful transmission of this data is at risk when used for
probing. Some applications may prefer to use a probe packet that
does not carry an application data block to avoid disruption to data
transfer.

### 6.1.4.  Initial Connectivity

An application that does not have other higher-layer information
confirming connectivity with the remote peer SHOULD implement a
connectivity mechanism using acknowledged probe packets before
entering the BASE state.

### 6.1.5.  Validating the Path

An application that does not have other higher-layer information
confirming correct delivery of datagrams SHOULD implement the
CONFIRMATION_TIMER to periodically send probe packets while in the
SEARCH_COMPLETE state.

### 6.1.6.  Handling of PTB Messages

An application that is able and wishes to receive PTB messages MUST
perform ICMP validation as specified in Section 5.2 of [RFC8085].
This requires that the application to check each received PTB
messages to validate it is received in response to transmitted
traffic and that the reported PTB_SIZE is less than the current
probed size (see Section 4.5.2). A validated PTB message MAY be used
as input to the DPLPMTUD algorithm, but MUST NOT be used directly to
set the PLPMTU.

### 6.2.  DPLPMTUD for SCTP

Section 10.2 of [RFC4821] specifies a recommended PLPMTUD probing
method for SCTP. It recommends the use of the PAD chunk, defined in
[RFC4820] to be attached to a minimum length HEARTBEAT chunk to
build a probe packet. This enables probing without affecting the
transfer of user messages and without interfering with congestion
control. This is preferred to using DATA chunks (with padding as
required) as path probes.

### 6.2.1.  SCTP/IPv4 and SCTP/IPv6

#### 6.2.1.1.  Initial Connectivity

The base protocol is specified in [RFC4960]. This provides an
acknowledged PL. A sender can therefore enter the BASE state as soon
as connectivity has been confirmed.

#### 6.2.1.2.  Sending SCTP Probe Packets

Probe packets consist of an SCTP common header followed by a
HEARTBEAT chunk and a PAD chunk. The PAD chunk is used to control
the length of the probe packet. The HEARTBEAT chunk is used to
trigger the sending of a HEARTBEAT ACK chunk. The reception of the
HEARTBEAT ACK chunk acknowledges reception of a successful probe.

The HEARTBEAT chunk carries a Heartbeat Information parameter which
should include, besides the information suggested in [RFC4960], the
probe size, which is the size of the complete datagram. The size of
the PAD chunk is therefore computed by reducing the probing size by
the IPv4 or IPv6 header size, the SCTP common header, the HEARTBEAT
request and the PAD chunk header. The payload of the PAD chunk
contains arbitrary data.

To avoid fragmentation of retransmitted data, probing starts right
after the PL handshake, before data is sent. Assuming this behavior
(i.e., the PMTU is smaller than or equal to the interface MTU), this
process will take a few round trip time periods depending on the
number of PMTU sizes probed. The Heartbeat timer can be used to
implement the PROBE_TIMER.

#### 6.2.1.3.  Validating the Path with SCTP

Since SCTP provides an acknowledged PL, a sender MUST NOT implement
the CONFIRMATION_TIMER while in the SEARCH_COMPLETE state.

#### 6.2.1.4.  PTB Message Handling by SCTP

Normal ICMP validation MUST be performed as specified in Appendix C
of [RFC4960]. This requires that the first 8 bytes of the SCTP
common header are quoted in the payload of the PTB message, which
can be the case for ICMPv4 and is normally the case for ICMPv6.

When a PTB message has been validated, the PTB_SIZE reported in the
PTB message SHOULD be used with the DPLPMTUD algorithm, providing
that the reported PTB_SIZE is less than the current probe size (see
Section 4.5).

### 6.2.2.  DPLPMTUD for SCTP/UDP

The UDP encapsulation of SCTP is specified in [RFC6951].

### 6.2.2.1.  Initial Connectivity

A sender can enter the BASE state as soon as SCTP connectivity has
been confirmed.

### 6.2.2.2.  Sending SCTP/UDP Probe Packets

Packet probing can be performed as specified in Section 6.2.1.2. The
maximum payload is reduced by 8 bytes, which has to be considered
when filling the PAD chunk.

### 6.2.2.3.  Validating the Path with SCTP/UDP

Since SCTP provides an acknowledged PL, a sender MUST NOT implement
the CONFIRMATION_TIMER while in the SEARCH_COMPLETE state.

### 6.2.2.4.  Handling of PTB Messages by SCTP/UDP

ICMP validation MUST be performed for PTB messages as specified in
Appendix C of [RFC4960]. This requires that the first 8 bytes of the
SCTP common header are contained in the PTB message, which can be
the case for ICMPv4 (but note the UDP header also consumes a part of
the quoted packet header) and is normally the case for ICMPv6. When
the validation is completed, the PTB_SIZE indicated in the PTB
message SHOULD be used with the DPLPMTUD providing that the reported
PTB_SIZE is less than the current probe size.

### 6.2.3.  DPLPMTUD for SCTP/DTLS

The Datagram Transport Layer Security (DTLS) encapsulation of SCTP
is specified in [RFC8261]. It is used for data channels in WebRTC
implementations.

### 6.2.3.1.  Initial Connectivity

A sender can enter the BASE state as soon as SCTP connectivity has
been confirmed.

### 6.2.3.2.  Sending SCTP/DTLS Probe Packets

Packet probing can be done as specified in Section 6.2.1.2.

### 6.2.3.3.  Validating the Path with SCTP/DTLS

Since SCTP provides an acknowledged PL, a sender MUST NOT implement
the CONFIRMATION_TIMER while in the SEARCH_COMPLETE state.

### 6.2.3.4.  Handling of PTB Messages by SCTP/DTLS

It is not possible to perform ICMP validation as specified in
[RFC4960], since even if the ICMP message payload contains
sufficient information, the reflected SCTP common header would be
encrypted. Therefore it is not possible to process PTB messages at
the PL.

### 6.3.  DPLPMTUD for QUIC

QUIC [I-D.ietf-quic-transport] is a UDP-based transport that
provides reception feedback. The UDP payload includes the QUIC
packet header, protected payload, and any authentication fields.
QUIC depends on a PMTU of at least 1280 bytes.

Section 14.1 of [I-D.ietf-quic-transport] describes the path
considerations when sending QUIC packets. It recommends the use of
PADDING frames to build the probe packet. Pure probe-only packets
are constructed with PADDING frames and PING frames to create a
padding only packet that will elicit an acknowledgement. Such
padding only packets enable probing without affecting the transfer
of other QUIC frames.

The recommendation for QUIC endpoints implementing DPLPMTUD is that
a MPS is maintained for each combination of local and remote IP
addresses [I-D.ietf-quic-transport]. If a QUIC endpoint determines
that the PMTU between any pair of local and remote IP addresses has
fallen below an acceptable MPS, it needs to immediately cease
sending QUIC packets on the affected path. This could result in
termination of the connection if an alternative path cannot be found
[I-D.ietf-quic-transport].

### 6.3.1.  Initial Connectivity

The base protocol is specified in [I-D.ietf-quic-transport]. This
provides an acknowledged PL. A sender can therefore enter the BASE
state as soon as connectivity has been confirmed.

### 6.3.2.  Sending QUIC Probe Packets

A probe packet consists of a QUIC Header and a payload containing
PADDING Frames and a PING Frame. PADDING Frames are a single octet
(0x00) and several of these can be used to create a probe packet of
size PROBED_SIZE. QUIC provides an acknowledged PL, a sender can
therefore enter the BASE state as soon as connectivity has been
confirmed.

The current specification of QUIC sets the following:

   *BASE_PMTU: 1200. A QUIC sender needs to pad initial packets to
    1200 bytes to confirm the path can support packets of a useful
    size.

   *MIN_PMTU: 1200 bytes. A QUIC sender that determines the PMTU has
    fallen below 1200 bytes MUST immediately stop sending on the
    affected path.

### 6.3.3.  Validating the Path with QUIC

QUIC provides an acknowledged PL. A sender therefore MUST NOT
implement the CONFIRMATION_TIMER while in the SEARCH_COMPLETE state.

### 6.3.4.  Handling of PTB Messages by QUIC

QUIC operates over the UDP transport, and the guidelines on ICMP
validation as specified in Section 5.2 of [RFC8085] therefore apply.
In addition to UDP Port validation QUIC can validate an ICMP message
by looking for valid Connection IDs in the quoted packet.

### 7.  Acknowledgements

### 8.  IANA Considerations

This memo includes no request to IANA.

If there are no requirements for IANA, the section will be removed
during conversion into an RFC by the RFC Editor.

### 9.  Security Considerations

The security considerations for the use of UDP and SCTP are provided
in the references RFCs. The interval between individual probe
packets MUST be at least one RTT, and the interval between rounds of
probing is determined by the PMTU_RAISE_TIMER.

A PL sender needs to ensure that the method used to confirm
reception of probe packets offers protection from off-path attackers
injecting packets into the path. This protection if provided in
IETF-defined protocols (e.g., TCP, SCTP) using a randomly-

initialized sequence number. A description of one way to do this
when using UDP is provided in section 5.1 of [RFC8085]).

There are cases where ICMP Packet Too Big (PTB) messages are not
delivered due to policy, configuration or equipment design (see
Section 1.1), this method therefore does not rely upon PTB messages
being received, but is able to utilize these when they are received
by the sender. PTB messages could potentially be used to cause a
node to inappropriately reduce the PLPMTU. A node supporting
DPLPMTUD MUST therefore appropriately validate the payload of PTB
messages to ensure these are received in response to transmitted
traffic (i.e., a reported error condition that corresponds to a
datagram actually sent by the path layer, see Section 4.5.1).

An on-path attacker, able to create a PTB message could forge PTB
messages that include a valid quoted IP packet. Such an attack could
be used to drive down the PLPMTU. There are two ways this method can
be mitigated against such attacks: First, by ensuring that a PL
sender never reduces the PLPMTU below the base size, solely in
response to receiving a PTB message. This is achieved by first
entering the BASE state when such a message is received. Second, the
design does not require processing of PTB messages, a PL sender
could therefore suspend processing of PTB messages (e.g., in a
robustness mode after detecting that subsequent probes actually
confirm that a size larger than the PTB_SIZE is supported by a
path).

Parallel forwarding paths SHOULD be considered. Section 5.4
identifies the need for robustness in the method when the path
information may be inconsistent.

A node performing DPLPMTUD could experience conflicting information
about the size of supported probe packets. This could occur when
there are multiple paths are concurrently in use and these exhibit a
different PMTU. If not considered, this could result in data being
black holed when the PLPMTU is larger than the smallest PMTU across
the current paths.

## 10.  References

### 10.1.  Normative References

[I-D.ietf-quic-transport]
          Iyengar, J. and M. Thomson, "QUIC: A UDP-Based
          Multiplexed and Secure Transport", Work in Progress,
          Internet-Draft, draft-ietf-quic-transport-20, 23 April

2019, <http://www.ietf.org/internet-drafts/draft-ietf-quic-transport-20.txt>.

[RFC0768]  Postel, J., "User Datagram Protocol", STD 6, RFC 768, DOI
           10.17487/RFC0768, August 1980, <https://www.rfc-editor.org/info/rfc768>.

[RFC0791]  Postel, J., "Internet Protocol", STD 5, RFC 791, DOI
           10.17487/RFC0791, September 1981, <https://www.rfc-editor.org/info/rfc791>.

[RFC1191]  Mogul, J.C. and S.E. Deering, "Path MTU discovery", RFC
           1191, DOI 10.17487/RFC1191, November 1990, <https://www.rfc-editor.org/info/rfc1191>.

[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>.

[RFC3828]  Larzon, L-A., Degermark, M., Pink, S., Jonsson, L-E.,
           Ed., and G. Fairhurst, Ed., "The Lightweight User
           Datagram Protocol (UDP-Lite)", RFC 3828, DOI 10.17487/RFC3828, July 2004, <https://www.rfc-editor.org/info/rfc3828>.

[RFC4820]  Tuexen, M., Stewart, R., and P. Lei, "Padding Chunk and
           Parameter for the Stream Control Transmission Protocol
           (SCTP)", RFC 4820, DOI 10.17487/RFC4820, March 2007,
           <https://www.rfc-editor.org/info/rfc4820>.

[RFC4960]  Stewart, R., Ed., "Stream Control Transmission Protocol",
           RFC 4960, DOI 10.17487/RFC4960, September 2007, <https://www.rfc-editor.org/info/rfc4960>.

[RFC6951]  Tuexen, M. and R. Stewart, "UDP Encapsulation of Stream
           Control Transmission Protocol (SCTP) Packets for End-Host
           to End-Host Communication", RFC 6951, DOI 10.17487/

              RFC6951, May 2013, <https://www.rfc-editor.org/info/
              rfc6951>.

   [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>.

   [RFC8174]  Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC
              2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174,
              May 2017, <https://www.rfc-editor.org/info/rfc8174>.

   [RFC8200]  Deering, S. and R. Hinden, "Internet Protocol, Version 6
              (IPv6) Specification", STD 86, RFC 8200, DOI 10.17487/
              RFC8200, July 2017, <https://www.rfc-editor.org/info/
              rfc8200>.

   [RFC8201]  McCann, J., Deering, S., Mogul, J., and R. Hinden, Ed.,
              "Path MTU Discovery for IP version 6", STD 87, RFC 8201,
              DOI 10.17487/RFC8201, July 2017, <https://www.rfc-
              editor.org/info/rfc8201>.

   [RFC8261]  Tuexen, M., Stewart, R., Jesup, R., and S. Loreto,
              "Datagram Transport Layer Security (DTLS) Encapsulation
              of SCTP Packets", RFC 8261, DOI 10.17487/RFC8261,
              November 2017, <https://www.rfc-editor.org/info/rfc8261>.

10.2.  Informative References

   [I-D.ietf-intarea-tunnels]  Touch, J. and M. Townsley, "IP Tunnels in
              the Internet Architecture", Work in Progress, Internet-
              Draft, draft-ietf-intarea-tunnels-10, 12 September 2019,
              <http://www.ietf.org/internet-drafts/draft-ietf-intarea-
              tunnels-10.txt>.

   [RFC0792]  Postel, J., "Internet Control Message Protocol", STD 5,
              RFC 792, DOI 10.17487/RFC0792, September 1981, <https://
              www.rfc-editor.org/info/rfc792>.

   [RFC1122]  Braden, R., Ed., "Requirements for Internet Hosts -
              Communication Layers", STD 3, RFC 1122, DOI 10.17487/
              RFC1122, October 1989, <https://www.rfc-editor.org/info/
              rfc1122>.

   [RFC1812]  Baker, F., Ed., "Requirements for IP Version 4 Routers",
              RFC 1812, DOI 10.17487/RFC1812, June 1995, <https://
              www.rfc-editor.org/info/rfc1812>.

   [RFC2923]  Lahey, K., "TCP Problems with Path MTU Discovery", RFC
              2923, DOI 10.17487/RFC2923, September 2000, <https://
              www.rfc-editor.org/info/rfc2923>.

**[RFC4340]**
Kohler, E., Handley, M., and S. Floyd, "Datagram Congestion Control Protocol (DCCP)", RFC 4340, DOI 10.17487/RFC4340, March 2006, <https://www.rfc-editor.org/info/rfc4340>.

**[RFC4443]** Conta, A., Deering, S., and M. Gupta, Ed., "Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification", STD 89, RFC 4443, DOI 10.17487/RFC4443, March 2006, <https://www.rfc-editor.org/info/rfc4443>.

**[RFC4821]** Mathis, M. and J. Heffner, "Packetization Layer Path MTU Discovery", RFC 4821, DOI 10.17487/RFC4821, March 2007, <https://www.rfc-editor.org/info/rfc4821>.

**[RFC4890]** Davies, E. and J. Mohacsi, "Recommendations for Filtering ICMPv6 Messages in Firewalls", RFC 4890, DOI 10.17487/ RFC4890, May 2007, <https://www.rfc-editor.org/info/ rfc4890>.

**Appendix A.  Revision Notes**

Note to RFC-Editor: please remove this entire section prior to publication.

Individual draft -00:

  *Comments and corrections are welcome directly to the authors or via the IETF TSVWG working group mailing list.

  *This update is proposed for WG comments.

Individual draft -01:

  *Contains the first representation of the algorithm, showing the states and timers

  *This update is proposed for WG comments.

Individual draft -02:

  *Contains updated representation of the algorithm, and textual corrections.

  *The text describing when to set the effective PMTU has not yet been validated by the authors

  *To determine security to off-path-attacks: We need to decide whether a received PTB message SHOULD/MUST be validated? The text

on how to handle a PTB message indicating a link MTU larger than
the probe has yet not been validated by the authors

   *No text currently describes how to handle inconsistent results
    from arbitrary re-routing along different parallel paths

   *This update is proposed for WG comments.

Working Group draft -00:

   *This draft follows a successful adoption call for TSVWG

   *There is still work to complete, please comment on this draft.

Working Group draft -01:

   *This draft includes improved introduction.

   *The draft is updated to require ICMP validation prior to
    accepting PTB messages - this to be confirmed by WG

   *Section added to discuss Selection of Probe Size - methods to be
    evaluated and recommendations to be considered

   *Section added to align with work proposed in the QUIC WG.

Working Group draft -02:

   *The draft was updated based on feedback from the WG, and a
    detailed review by Magnus Westerlund.

   *The document updates RFC 4821.

   *Requirements list updated.

   *Added more explicit discussion of a simpler black-hole detection
    mode.

   *This draft includes reorganisation of the section on IETF
    protocols.

   *Added more discussion of implementation within an application.

   *Added text on flapping paths.

   *Replaced 'effective MTU' with new term PLPMTU.

Working Group draft -03:

   *Updated figures

*Added more discussion on blackhole detection

   *Added figure describing just blackhole detection

   *Added figure relating MPS sizes

Working Group draft -04:

  *Described phases and named these consistently.

  *Corrected transition from confirmation directly to the search
   phase (Base has been checked).

  *Redrawn state diagrams.

  *Renamed BASE_MTU to BASE_PMTU (because it is a base for the
   PMTU).

  *Clarified Error state.

  *Clarified suspending DPLPMTUD.

  *Verified normative text in requirements section.

  *Removed duplicate text.

  *Changed all text to refer to /packet probe/probe packet/ /
   validation/verification/ added term /Probe Confirmation/ and
   clarified BlackHole detection.

Working Group draft -05:

  *Updated security considerations.

  *Feedback after speaking with Joe Touch helped improve UDP-Options
   description.

Working Group draft -06:

  *Updated description of ICMP issues in section 1.1

  *Update to description of QUIC.

Working group draft -07:

  *Moved description of the PTB processing method from the PTB
   requirements section.

  *Clarified what is performed in the PTB validation check.

*Updated security consideration to explain PTB security without
    needing to read the rest of the document.

   *Reformatted state machine diagram

Working group draft -08:

   *Moved to rfcxml v3+

   *Rendered diagrams to svg in html version.

   *Removed Appendix A. Event-driven state changes.

   *Removed section on DPLPMTUD with UDP Options.

   *Shortened the description of phases.

Working group draft -09:

   *Remove final mention of UDP Options

   *Add Initial Connectivity sections to each PL

   *Add to disable outgoing pmtu enforcement of packets

Working group draft -10:

   *Address comments from Lars Eggert

   *Reinforce that PROBE_COUNT is successive attempts to probe for
    any size

   *Redefine MAx_PROBES to 3

   *Address PTB_SIZE of 0 or less that MIN_PMTU

Working group draft -11:

   *Restore a sentence removed in previous rev

   *De-acronymise QUIC

   *Address some nits

Working group draft -12:

   *Add TSVWG, QUIC and implementers to acknowledgements

   *Shorten a diagram line

   *Address nits from Julius and Wes

*Be clearer when talking about IP layer caches

**Authors' Addresses**

   Godred Fairhurst
   University of Aberdeen
   School of Engineering, Fraser Noble Building
   Aberdeen
   AB24 3UE
   United Kingdom


   Email: gorry@erg.abdn.ac.uk


   Tom Jones
   University of Aberdeen
   School of Engineering, Fraser Noble Building
   Aberdeen
   AB24 3UE
   United Kingdom


   Email: tom@erg.abdn.ac.uk


   Michael Tuexen
   Muenster University of Applied Sciences
   Stegerwaldstrasse 39
   48565 Steinfurt
   Germany


   Email: tuexen@fh-muenster.de


   Irene Ruengeler
   Muenster University of Applied Sciences
   Stegerwaldstrasse 39
   48565 Steinfurt
   Germany


   Email: i.ruengeler@fh-muenster.de


   Timo Voelker
   Muenster University of Applied Sciences
   Stegerwaldstrasse 39
   48565 Steinfurt
   Germany


   Email: timo.voelker@fh-muenster.de