

IPPM  
Internet-Draft  
Intended status: Informational  
Expires: April 28, 2022

M. Cociglio  
Telecom Italia - TIM  
A. Ferrieux  
Orange Labs  
G. Fioccola  
Huawei Technologies  
I. Lubashev  
Akamai Technologies  
F. Bulgarella  
Telecom Italia - TIM  
I. Hamchaoui  
Orange Labs  
M. Nilo  
Telecom Italia - TIM  
R. Sisto  
Politecnico di Torino  
D. Tikhonov  
LiteSpeed Technologies  
October 25, 2021

Explicit Flow Measurements Techniques  
draft-ietf-ippm-explicit-flow-measurements-00

Abstract

This document describes protocol independent methods called Explicit Flow Measurement Techniques that employ few marking bits, inside the header of each packet, for loss and delay measurement. The endpoints, marking the traffic, signal these metrics to intermediate observers allowing them to measure connection performance, and to locate the network segment where impairments happen. Different alternatives are considered within this document. These signaling methods apply to all protocols but they are especially valuable when applied to protocols that encrypt transport header and do not allow traditional methods for delay and loss detection.

Status of This Memo

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

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

Internet-Draft

Delay and Loss bits

October 2021

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

This Internet-Draft will expire on April 28, 2022.

## Copyright Notice

Copyright (c) 2021 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](https://trustee.ietf.org/bcp78) and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

<a href="#">1.</a>	Introduction . . . . .	<a href="#">3</a>
<a href="#">2.</a>	Notational Conventions . . . . .	<a href="#">5</a>
<a href="#">3.</a>	Latency Bits . . . . .	<a href="#">5</a>
<a href="#">3.1.</a>	Spin Bit . . . . .	<a href="#">5</a>
<a href="#">3.2.</a>	Delay Bit . . . . .	<a href="#">6</a>
<a href="#">3.2.1.</a>	Generation Phase . . . . .	<a href="#">9</a>
<a href="#">3.2.2.</a>	Reflection Phase . . . . .	<a href="#">9</a>
<a href="#">3.2.3.</a>	T_Max Selection . . . . .	<a href="#">10</a>
<a href="#">3.2.4.</a>	Delay Measurement using Delay Bit . . . . .	<a href="#">11</a>
<a href="#">3.2.5.</a>	Observer's Algorithm . . . . .	<a href="#">13</a>
<a href="#">3.2.6.</a>	Two Bits Delay Measurement: Spin Bit + Delay Bit . . . . .	<a href="#">14</a>
<a href="#">3.2.7.</a>	Hidden Delay Bit - Delay Bit with Privacy Protection . . . . .	14
<a href="#">4.</a>	Loss Bits . . . . .	<a href="#">14</a>
<a href="#">4.1.</a>	T Bit - Round Trip Loss Bit . . . . .	<a href="#">15</a>
<a href="#">4.1.1.</a>	Round Trip Packet Loss Measurement . . . . .	<a href="#">16</a>
<a href="#">4.1.2.</a>	Setting the Round Trip Loss Bit on Outgoing Packets . . . . .	18
<a href="#">4.1.3.</a>	Observer's Logic for Round Trip Loss Signal . . . . .	<a href="#">19</a>
<a href="#">4.1.4.</a>	Loss Coverage and Signal Timing . . . . .	<a href="#">20</a>
<a href="#">4.2.</a>	Q Bit - Square Bit . . . . .	<a href="#">20</a>

<a href="#">4.2.1.</a>	Q Block Length Selection . . . . .	<a href="#">20</a>
<a href="#">4.2.2.</a>	Upstream Loss . . . . .	<a href="#">21</a>
<a href="#">4.2.3.</a>	Identifying Q Block Boundaries . . . . .	<a href="#">22</a>
<a href="#">4.3.</a>	L Bit - Loss Event Bit . . . . .	<a href="#">22</a>
<a href="#">4.3.1.</a>	End-To-End Loss . . . . .	<a href="#">23</a>

<a href="#">4.3.2.</a>	Loss Profile Characterization . . . . .	<a href="#">23</a>
<a href="#">4.4.</a>	L+Q Bits - Upstream, Downstream, and End-to-End Loss Measurements . . . . .	<a href="#">23</a>
<a href="#">4.4.1.</a>	Correlating End-to-End and Upstream Loss . . . . .	<a href="#">24</a>
<a href="#">4.5.</a>	R Bit - Reflection Square Bit . . . . .	<a href="#">25</a>
<a href="#">4.5.1.</a>	R+Q Bits - Using R and Q Bits for Passive Loss Measurement . . . . .	<a href="#">26</a>
<a href="#">4.5.2.</a>	Enhancement of R Block Length Computation . . . . .	<a href="#">30</a>
<a href="#">4.5.3.</a>	Improved Resilience to Packet Reordering . . . . .	<a href="#">30</a>
<a href="#">4.6.</a>	Improved Q and R Bits Resilience to Burst Losses . . . . .	<a href="#">30</a>
<a href="#">5.</a>	Summary of Delay and Loss Marking Methods . . . . .	<a href="#">31</a>
<a href="#">6.</a>	ECN-Echo Event Bit . . . . .	<a href="#">33</a>
<a href="#">6.1.</a>	Setting the ECN-Echo Event Bit on Outgoing Packets . . . . .	<a href="#">33</a>
<a href="#">6.2.</a>	Using E Bit for Passive ECN-Reported Congestion Measurement . . . . .	<a href="#">33</a>
<a href="#">7.</a>	Protocol Ossification Considerations . . . . .	<a href="#">34</a>
<a href="#">8.</a>	Examples of Application . . . . .	<a href="#">34</a>
<a href="#">8.1.</a>	QUIC . . . . .	<a href="#">34</a>
<a href="#">8.2.</a>	TCP . . . . .	<a href="#">35</a>
<a href="#">9.</a>	Security Considerations . . . . .	<a href="#">35</a>
<a href="#">9.1.</a>	Optimistic ACK Attack . . . . .	<a href="#">36</a>
<a href="#">10.</a>	Privacy Considerations . . . . .	<a href="#">36</a>
<a href="#">11.</a>	IANA Considerations . . . . .	<a href="#">37</a>
<a href="#">12.</a>	Change Log . . . . .	<a href="#">37</a>
<a href="#">13.</a>	Contributors . . . . .	<a href="#">37</a>
<a href="#">14.</a>	Acknowledgements . . . . .	<a href="#">37</a>
<a href="#">15.</a>	References . . . . .	<a href="#">37</a>
<a href="#">15.1.</a>	Normative References . . . . .	<a href="#">37</a>
<a href="#">15.2.</a>	Informative References . . . . .	<a href="#">38</a>
	Authors' Addresses . . . . .	<a href="#">40</a>

[1.](#) Introduction

Packet loss and delay are hard and pervasive problems of day-to-day network operation. Proactively detecting, measuring, and locating them is crucial to maintaining high QoS and timely resolution of

crippling end-to-end throughput issues. To this effect, in a TCP-dominated world, network operators have been heavily relying on information present in the clear in TCP headers: sequence and acknowledgment numbers and SACKs when enabled (see [\[RFC8517\]](#)). These allow for quantitative estimation of packet loss and delay by passive on-path observation. Additionally, the problem can be quickly identified in the network path by moving the passive observer around.

With encrypted protocols, the equivalent transport headers are encrypted and passive packet loss and delay observations are not possible, as described in [\[RFC9065\]](#).

Measuring TCP loss and delay between similar endpoints cannot be relied upon to evaluate encrypted protocol loss and delay. Different protocols could be routed by the network differently, and the fraction of Internet traffic delivered using protocols other than TCP is increasing every year. It is imperative to measure packet loss and delay experienced by encrypted protocol users directly.

This document defines Explicit Flow Measurement Techniques. These hybrid measurement path signals (see [\[IPM-Methods\]](#)) are to be embedded into a transport layer protocol and are explicitly intended for exposing RTT and loss rate information to on-path measurement devices. They are designed to facilitate network operations and management and are "beneficial" for maintaining the quality of service (see [\[RFC9065\]](#)). These measurement mechanisms are applicable to any transport-layer protocol, and, as an example, the document describes QUIC and TCP bindings.

The Explicit Flow Measurement Techniques described in this document can be used alone or in combination with other Explicit Flow Measurement Techniques. Each technique uses a small number of bits and exposes a specific measurement.

Following the recommendation in [\[RFC8558\]](#) of making path signals explicit, this document proposes adding a small number of dedicated measurement bits to the clear portion of the protocol headers. These bits can be added to an unencrypted portion of a header belonging to any protocol layer, e.g. IP (see [\[IP\]](#)) and IPv6 (see [\[IPv6\]](#)) headers or extensions, such as [\[IPv6AltMark\]](#), UDP surplus space (see [\[UDP-OPTIONS\]](#) and [\[UDP-SURPLUS\]](#)), reserved bits in a QUIC v1 header,

as already done with the latency spin bit (see [[QUIC-TRANSPORT](#)]).

The measurements are not designed for use in automated control of the network in environments where signal bits are set by untrusted hosts. Instead, the signal is to be used for troubleshooting individual flows as well as for monitoring the network by aggregating information from multiple flows and raising operator alarms if aggregate statistics indicate a potential problem.

The spin bit, delay bit and loss bits explained in this document are inspired by [[AltMark](#)], [[SPIN-BIT](#)], [[I-D.trammell-tsvwg-spin](#)] and [[I-D.trammell-ippm-spin](#)].

Additional details about the Performance Measurements for QUIC are described in the paper [[ANRW19-PM-QUIC](#)].

## [2.](#) Notational Conventions

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

## [3.](#) Latency Bits

This section introduces bits that can be used for round trip latency measurements. Whenever this section of the specification refers to packets, it is referring only to packets with protocol headers that include the latency bits.

[[QUIC-TRANSPORT](#)] introduces an explicit per-flow transport-layer signal for hybrid measurement of RTT. This signal consists of a spin bit that toggles once per RTT. [[SPIN-BIT](#)] discusses an additional two-bit Valid Edge Counter (VEC) to compensate for loss and reordering of the spin bit and increase fidelity of the signal in less than ideal network conditions.

This document introduces a stand-alone single-bit delay signal that can be used by passive observers to measure the RTT of a network

flow, avoiding the spin bit ambiguities that arise as soon as network conditions deteriorate.

### [3.1.](#) Spin Bit

This section is a small recap of the spin bit working mechanism. For a comprehensive explanation of the algorithm, please see [[SPIN-BIT](#)].

The spin bit is an alternate marking [[AltMark](#)] generated signal, where the size of the alternation changes with the flight size each RTT.

The latency spin bit is a single bit signal that toggles once per RTT, enabling latency monitoring of a connection-oriented communication from intermediate observation points.

A "spin period" is a set of packets with the same spin bit value sent during one RTT time interval. A "spin period value" is the value of the spin bit shared by all packets in a spin period.

The client and server maintain an internal per-connection spin value (i.e. 0 or 1) used to set the spin bit on outgoing packets. Both endpoints initialize the spin value to 0 when a new connection starts. Then:

- when the client receives a packet with the packet number larger than any number seen so far, it sets the connection spin value to the opposite value contained in the received packet;
- when the server receives a packet with the packet number larger than any number seen so far, it sets the connection spin value to the same value contained in the received packet.

The computed spin value is used by the endpoints for setting the spin bit on outgoing packets. This mechanism allows the endpoints to generate a square wave such that, by measuring the distance in time between pairs of consecutive edges observed in the same direction, a passive on-path observer can compute the round trip delay of that network flow.

Spin bit enables round trip latency measurement by observing a single direction of the traffic flow.

Note that packet reordering can cause spurious edges that require heuristics to correct. The spin bit performance deteriorates as soon as network impairments arise as explained in [Section 3.2](#).

### [3.2](#). Delay Bit

The delay bit has been designed to overcome accuracy limitations experienced by the spin bit under difficult network conditions:

- packet reordering leads to generation of spurious edges and errors in delay estimation;
- loss of edges causes wrong estimation of spin periods and therefore wrong RTT measurements;
- application-limited senders cause the spin bit to measure the application delays instead of network delays.

Unlike the spin bit, which is set in every packet transmitted on the network, the delay bit is set only once per round trip.

When the delay bit is used, a single packet with a marked bit (the delay bit) bounces between a client and a server during the entire connection lifetime. This single packet is called "delay sample".

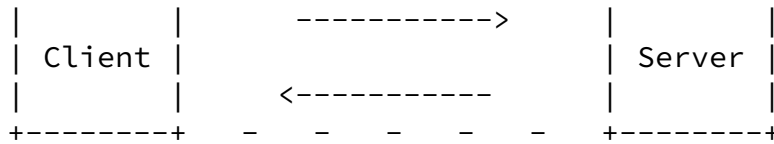
An observer placed at an intermediate point, observing a single direction of traffic, tracking the delay sample and the relative timestamp, can measure the round trip delay of the connection.

The delay sample lifetime is comprised of two phases: initialization and reflection. The initialization is the generation of the delay sample, while the reflection realizes the bounce behavior of this single packet between the two endpoints.

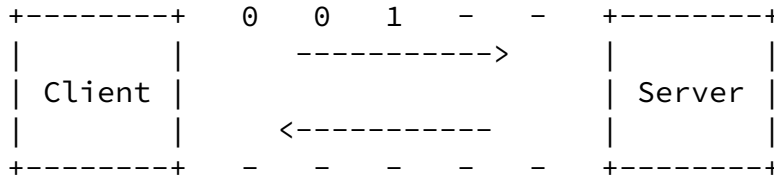
The next figure describes the elementary Delay bit mechanism.

+-----+ - - - - - +-----+

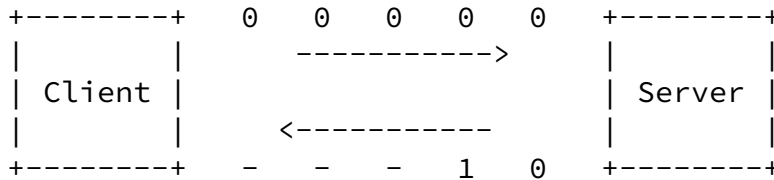




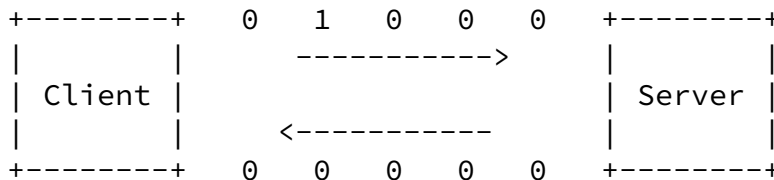
(a) No traffic at beginning.



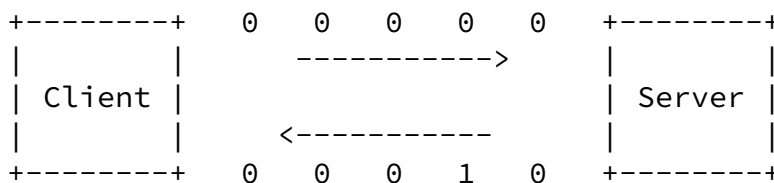
(b) The Client starts sending data and sets the first packet as Delay Sample.



(c) The Server starts sending data and reflects the Delay Sample.



(d) The Client reflects the Delay Sample.



(e) The Server reflects the Delay Sample and so on.

### Delay bit mechanism

### [3.2.1.](#) Generation Phase

Only client is actively involved in the generation phase. It maintains an internal per-flow timestamp variable ("ds\_time") updated every time a delay sample is transmitted.

When connection starts, the client generates a new delay sample initializing the delay bit of the first outgoing packet to 1. Then it updates the "ds\_time" variable with the timestamp of its transmission.

The server initializes the delay bit to 0 at the beginning of the connection, and its only task during the connection is described in [Section 3.2.2.](#)

In absence of network impairments, the delay sample should bounce between client and server continuously, for the entire duration of the connection. That is highly unlikely for two reasons:

1. the packet carrying the delay bit might be lost;
2. an endpoint could stop or delay sending packets because the application is limiting the amount of traffic transmitted;

To deal with these problems, the client generates a new delay sample if more than a predetermined time ("T\_Max") has elapsed since the last delay sample transmission (including reflections). Note that "T\_Max" should be greater than the max measurable RTT on the network. See [Section 3.2.3](#) for details.

### [3.2.2.](#) Reflection Phase

Reflection is the process that enables the bouncing of the delay sample between a client and a server. The behavior of the two endpoints is almost the same.

- Server side reflection: when a delay sample arrives, the server marks the first packet in the opposite direction as the delay sample.
- Client side reflection: when a delay sample arrives, the client marks the first packet in the opposite direction as the delay sample. It also updates the "ds\_time" variable when the outgoing delay sample is actually forwarded.

In both cases, if the outgoing delay sample is being transmitted with

a delay greater than a predetermined threshold after the reception of

the incoming delay sample (1ms by default), the delay sample is not reflected, and the outgoing delay bit is kept at 0.

By doing so, the algorithm can reject measurements that would overestimate the delay due to lack of traffic on the endpoints. Hence, the maximum estimation error would amount to twice the threshold (e.g. 2ms) per measurement.

### [3.2.3.](#) T\_Max Selection

The internal "ds\_time" variable allows a client to identify delay sample losses. Considering that a lost delay sample is regenerated at the end of an explicit time ("T\_Max") since the last generation, this same value can be used by an observer to reject a measure and start a new one.

In other words, if the difference in time between two delay samples is greater or equal than "T\_Max", then these cannot be used to produce a delay measure. Therefore the value of "T\_Max" must also be known to the on-path network probes.

There are two alternatives to select the "T\_Max" value so that both client and observers know it. The first one requires that "T\_Max" is known a priori ("T\_Max\_p") and therefore set within the protocol specifications that implements the marking mechanism (e.g. 1 second which usually is greater than the max expectable RTT). The second alternative requires a dynamic mechanism able to adapt the duration of the "T\_Max" to the delay of the connection ("T\_Max\_c").

For instance, client and observers could use the connection RTT as a basis for calculating an effective "T\_Max". They should use a predetermined initial value so that "T\_Max = T\_Max\_p" (e.g. 1 second) and then, when a valid RTT is measured, change "T\_Max" accordingly so that "T\_Max = T\_Max\_c". In any case, the selected "T\_Max" should be large enough to absorb any possible variations in the connection delay.

"T\_Max\_c" could be computed as two times the measured "RTT" plus a fixed amount of time ("100ms") to prevent low "T\_Max" values in case

of very small RTTs. The resulting formula is: "T\_Max\_c = 2RTT + 100ms". If "T\_Max\_c" is greater than "T\_Max\_p" then "T\_Max\_c" is forced to "T\_Max\_p" value.

Note that the observer's "T\_Max" should always be less than or equal to the client's "T\_Max" to avoid considering as a valid measurement what is actually the client's "T\_Max". To obtain this result, the client waits for two consecutive incoming samples and computes the two related RTTs. Then it takes the largest of them as the basis of

the "T\_Max\_c" formula. At this point, observers have already measured a valid RTT and then computed their "T\_Max\_c".

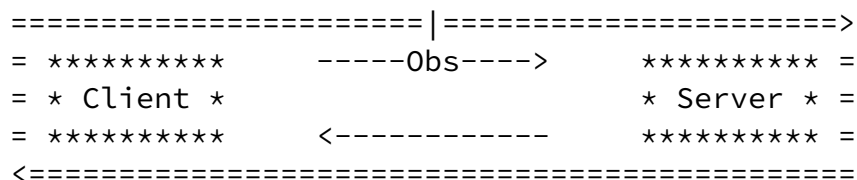
### [3.2.4.](#) Delay Measurement using Delay Bit

When the Delay Bit is used, a passive observer can use delay samples directly and avoid inherent ambiguities in the calculation of the RTT as can be seen in spin bit analysis.

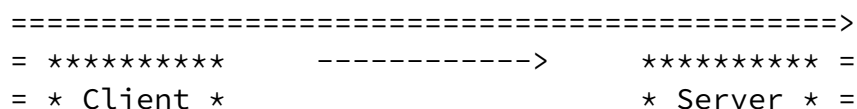
#### [3.2.4.1.](#) RTT Measurement

The delay sample generation process ensures that only one packet marked with the delay bit set to 1 runs back and forth between two endpoints per round trip time. To determine the RTT measurement of a flow, an on-path passive observer computes the time difference between two delay samples observed in a single direction.

To ensure a valid measurement, the observer must verify that the distance in time between the two samples taken into account is less than "T\_Max".



(a) client-server RTT



```

= ***** <----Obs----- ***** =
<=====|=====

```

(b) server-client RTT

Round-trip time (both direction)

### 3.2.4.2. Half-RTT Measurement

An observer that is able to observe both forward and return traffic directions can use the delay samples to measure "upstream" and "downstream" RTT components, also known as the half-RTT measurements. It does this by measuring the time between a delay sample observed in one direction and the delay sample previously observed in the opposite direction.

As with RTT measurement, the observer must verify that the distance in time between the two samples taken into account is less than "T\_Max".

Note that upstream and downstream sections of paths between the endpoints and the observer, i.e. observer-to-client vs client-to-observer and observer-to-server vs server-to-observer, may have different delay characteristics due to the difference in network congestion and other factors.

```

=====>
= ***** -----|-----> *****
= * Client *          Obs          * Server *
= ***** <-----|----- *****
<=====

```

(a) client-observer half-RTT

```

=====>
***** -----|-----> ***** =
* Client *          Obs          * Server * =
***** <-----|----- ***** =
<=====

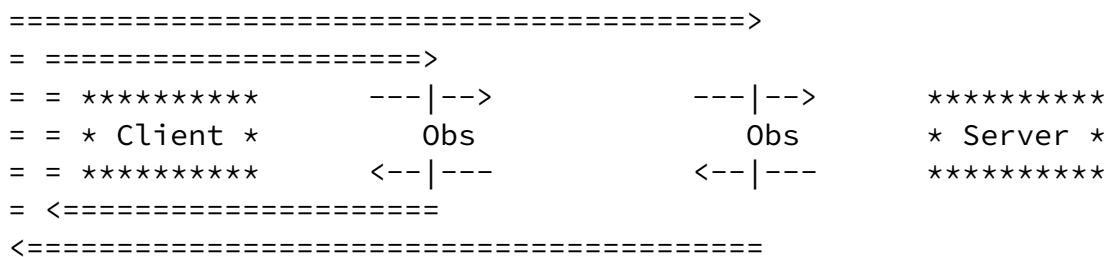
```

(b) observer-server half-RTT

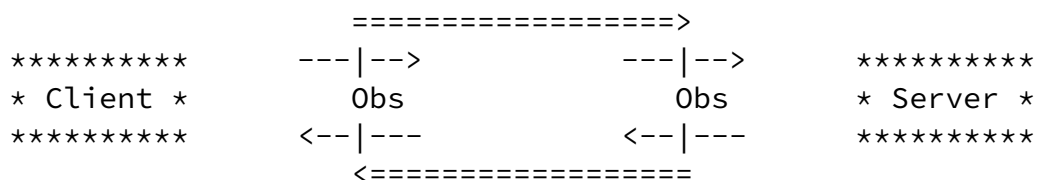
Half Round-trip time (both direction)

[3.2.4.3.](#) Intra-Domain RTT Measurement

Intra-domain RTT is the portion of the entire RTT used by a flow to traverse the network of a provider. To measure intra-domain RTT, two observers capable of observing traffic in both directions must be employed simultaneously at ingress and egress of the network to be measured. Intra-domain RTT is difference between the two computed upstream (or downstream) RTT components.



(a) client-observer RTT components (half-RTTs)



(b) the intra-domain RTT resulting from the subtraction of the above RTT components

Intra-domain Round-trip time (client-observer: upstream)

### [3.2.5.](#) Observer's Algorithm

An on-path observer maintains an internal per-flow variable to keep track of time at which the last delay sample has been observed.

A unidirectional observer, upon detecting a delay sample:

- if a delay sample was also detected previously in the same direction and the distance in time between them is less than " $T_{Max} - K$ ", then the two delay samples can be used to calculate RTT measurement. " $K$ " is a protection threshold to absorb differences in " $T_{Max}$ " computation and delay variations between two consecutive delay samples (e.g. " $K = 10\% T_{Max}$ ").

If the observer can observe both forward and return traffic flows, and it is able to determine which direction contains the client and the server (e.g. by observing the connection handshake), upon detecting a delay sample:

- if a delay sample was also detected in the opposite direction and the distance in time between them is less than " $T_{Max} - K$ ", then the two delay samples can be used to measure the observer-client half-RTT or the observer-server half-RTT, according to the direction of the last delay sample observed.

### [3.2.6.](#) Two Bits Delay Measurement: Spin Bit + Delay Bit

Spin and Delay bit algorithms work independently. If both marking methods are used in the same connection, observers can choose the best measurement between the two available:

- when a precise measurement can be produced using the delay bit, observers choose it;
- when a delay bit measurement is not available, observers choose the approximate spin bit one.

### [3.2.7.](#) Hidden Delay Bit - Delay Bit with Privacy Protection

Theoretically, delay measurements can be used to roughly evaluate the distance of the client from the server (using the RTT) or from any intermediate observer (using the client-observer half-RTT). To protect users privacy, the algorithm of the delay bit can be slightly modified to mask the RTT of the connection to an intermediate observer. This result can be achieved using a simple expedient which consists in delaying the client-side reflection of the delay sample by a predetermined time value. This would lead an intermediate observer to inevitably measure a delay greater than the real one.

The Additional Delay should be randomly selected by the client and kept constant for a certain amount of time across multiple connections. This ensures that the client-server jitter remains the same as if no Additional Delay had been inserted. For instance, a new Additional Delay value could be generated whenever the client's IP address changes.

Using this technique, despite the Additional Delay introduced, it is still possible to correctly measure the right component of RTT (observer-server) and all the intra-domain measurements used to distribute the delay in the network. Furthermore, differently from the Delay Bit, the hidden Delay Bit makes the use of the client reflection threshold (1ms) redundant. Removing this threshold leads to the further advantage of increasing the number of valid measurements produced by the algorithm.

## [4.](#) Loss Bits

This section introduces bits that can be used for loss measurements. Whenever this section of the specification refers to packets, it is referring only to packets with protocol headers that include the loss bits - the only packets whose loss can be measured.

- T: the "round Trip loss" bit is used in combination with the Spin bit to measure round-trip loss. See [Section 4.1](#).
- Q: the "square signal" bit is used to measure upstream loss. See



## [Section 4.2.](#)

- L: the "Loss event" bit is used to measure end-to-end loss. See [Section 4.3.](#)
- R: the "Reflection square signal" bit is used in combination with Q bit to measure end-to-end loss. See [Section 4.1.](#)

Loss measurements enabled by T, Q, and L bits can be implemented by those loss bits alone (T bit requires a working Spin Bit). Two-bit combinations Q+L and Q+R enable additional measurement opportunities discussed below.

Each endpoint maintains appropriate counters independently and separately for each separately identifiable flow (each sub-flow for multipath connections).

Since loss is reported independently for each flow, all bits (except for L bit) require a certain minimum number of packets to be exchanged per flow before any signal can be measured. Therefore, loss measurements work best for flows that transfer more than a minimal amount of data.

### [4.1.](#) T Bit - Round Trip Loss Bit

The round Trip loss bit is used to mark a variable number of packets exchanged twice between the endpoints realizing a two round-trip reflection. A passive on-path observer, observing either direction, can count and compare the number of marked packets seen during the two reflections, estimating the loss rate experienced by the connection. The overall exchange comprises:

- The client selects, generates and consequently transmits a first train of packets, by setting the T bit to 1;
- The server, upon receiving each packet included in the first train, reflects to the client a respective second train of packets of the same size as the first train received, by setting the T bit to 1;
- The client, upon receiving each packet included in the second train, reflects to the server a respective third train of packets of the same size as the second train received, by setting the T bit to 1;

- The server, upon receiving each packet included in the third train, finally reflects to the client a respective fourth train of packets of the same size as the third train received, by setting the T bit to 1.

Packets belonging to the first round trip (first and second train) represent the Generation Phase, while those belonging to the second round trip (third and fourth train) represent the Reflection Phase.

A passive on-path observer can count and compare the number of marked packets seen during the two round trips (i.e. the first and third or the second and the fourth trains of packets, depending on which direction is observed) and estimate the loss rate experienced by the connection. This process is repeated continuously to obtain more measurements as long as the endpoints exchange traffic. These measurements can be called Round Trip losses.

Since packet rates in two directions may be different, the number of marked packets in the train is determined by the direction with the lowest packet rate. See [Section 4.1.2](#) for details on packet generation and for a mechanism to allow an observer to distinguish between trains belonging to different phases (Generation and Reflection).

#### [4.1.1](#). Round Trip Packet Loss Measurement

Since the measurements are performed on a portion of the traffic exchanged between the client and the server, the observer calculates the end-to-end Round Trip Packet Loss (RTPL) that, statistically, will correspond to the loss rate experienced by the connection along the entire network path.

Internet-Draft

Delay and Loss bits

October 2021

```

=====|=====>
= ***** -----Obs----> ***** =
= * Client *                               * Server * =
= ***** <-----< ***** =
<=====|=====

```

(a) client-server RTPL

```

=====|=====>
= ***** -----> ***** =
= * Client *                               * Server * =
= ***** <-----Obs----- ***** =
<=====|=====

```

(b) server-client RTPL

Round-trip packet loss (both direction)

This methodology also allows the Half-RTPL measurement and the Intra-domain RTPL measurement in a way similar to RTT measurement.

```

=====|=====>
= ***** -----|-----> *****
= * Client *           Obs           * Server *
= ***** <-----|-----< *****
<=====|=====

```

(a) client-observer half-RTPL

```

=====|=====>
***** -----|-----> ***** =
* Client *           Obs           * Server * =
***** <-----|-----< ***** =
<=====|=====

```

(b) observer-server half-RTPL

Half Round-trip packet loss (both direction)

Internet-Draft

Delay and Loss bits

October 2021

```

=====>
                                     =====> =
*****      ---|-->          ---|-->      ***** = =
* Client *      Obs          Obs          * Server * = =
*****      <--|---          <--|---          ***** = =
                                     <===== =
<=====

```

(a) observer-server RTPL components (half-RTPLs)

```

=====>
*****      ---|-->          ---|-->      *****
* Client *      Obs          Obs          * Server *
*****      <--|---          <--|---          *****
<=====

```

(b) the intra-domain RTPL resulting from the subtraction of the above RTPL components

Intra-domain Round-trip packet loss (observer-server)

#### [4.1.2.](#) Setting the Round Trip Loss Bit on Outgoing Packets

The round Trip loss signal requires a working Spin-bit signal to separate trains of marked packets (packets with T bit set to 1). A "pause" of at least one empty spin-bit period between each phase of the algorithm serves as such separator for the on-path observer.

The client is in charge of launching trains of marked packets and does so according to the algorithm:

1. Generation Phase. The client starts generating marked packets for two consecutive spin-bit periods; it maintains a "generation

token" count that is reset to zero at the beginning of the algorithm phase and is incremented every time a packet arrives. When the client transmits a packet and a "generation token" is available, the client marks the packet and retires a "generation token". If no token is available, the outgoing packet is transmitted unmarked. At the end of the first spin-bit period spent in generation, the reflection counter is unlocked to start counting incoming marked packets that will be reflected later;

2. Pause Phase. When the generation is completed, the client pauses till it has observed one entire spin bit period with no marked packets. That spin bit period is used by the observer as a separator between generated and reflected packets. During this marking pause, all the outgoing packets are transmitted with T

bit set to 0. The reflection counter is still incremented every time a marked packet arrives;

3. Reflection Phase. The client starts transmitting marked packets, decrementing the reflection counter for each transmitted marked packet until the reflection counter reached zero. The "generation token" method from the generation phase is used during this phase as well. At the end of the first spin-period spent in reflection, the reflection counter is locked to avoid incoming reflected packets incrementing it;
4. Pause Phase 2. The pause phase is repeated after the reflection phase and serves as a separator between the reflected packet train and a new packet train.

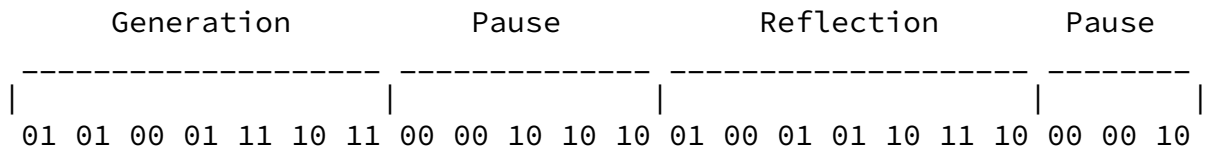
The generation token counter should be capped to limit the effects of a subsequent sudden reduction in the other endpoint's packet rate that could prevent that endpoint from reflecting collected packets. The most conservative cap value is "1".

A server maintains a "marking counter" that starts at zero and is incremented every time a marked packet arrives. When the server transmits a packet and the "marking counter" is positive, the server marks the packet and decrements the "marking counter". If the "marking counter" is zero, the outgoing packet is transmitted unmarked.

### 4.1.3. Observer's Logic for Round Trip Loss Signal

The on-path observer counts marked packets and separates different trains by detecting spin-bit periods (at least one) with no marked packets. The Round Trip Packet Loss (RTPL) is the difference between the size of the Generation train and the Reflection train.

In the following example, packets are represented by two bits (first one is the spin bit, second one is the loss bit):



Round Trip Loss signal example

Note that 5 marked packets have been generated of which 4 have been reflected.

### 4.1.4. Loss Coverage and Signal Timing

A cycle of the round Trip loss signaling algorithm contains 2 RTTs of Generation phase, 2 RTTs of Reflection phase, and two Pause phases at least 1 RTT in duration each. Hence, the loss signal is delayed by about 6 RTTs since the loss events.

The observer can only detect loss of marked packets that occurs after its initial observation of the Generation phase and before its subsequent observation of the Reflection phase. Hence, if the loss occurs on the path that sends packets at a lower rate (typically ACKs in such asymmetric scenarios), "2/6" ("1/3") of the packets will be sampled for loss detection.

If the loss occurs on the path that sends packets at a higher rate, " $\text{lowPacketRate}/(3*\text{highPacketRate})$ " of the packets will be sampled for loss detection. For protocols that use ACKs, the portion of packets sampled for loss in the higher rate direction during unidirectional data transfer is " $1/(3*\text{packetsPerAck})$ ", where the value of

"packetsPerAck" can vary by protocol, by implementation, and by network conditions.

## [4.2.](#) Q Bit - Square Bit

The square bit (Q bit) takes its name from the square wave generated by its signal. Every outgoing packet contains the Q bit value, which is initialized to the 0 and inverted after sending N packets (a square Block or simply Q Block). Hence, Q Period is  $2*N$ . The Q bit represents "packet color" as defined by [\[AltMark\]](#).

Observation points can estimate upstream losses by watching a single direction of the traffic flow and counting the number of packets in each observed Q Block, as described in [Section 4.2.2](#).

### [4.2.1.](#) Q Block Length Selection

The length of the block must be known to the on-path network probes. There are two alternatives to selecting the Q Block length. The first one requires that the length is known a priori and therefore set within the protocol specifications that implements the marking mechanism. The second requires the sender to select it.

In this latter scenario, the sender is expected to choose N (Q Block length) based on the expected amount of loss and reordering on the path. The choice of N strikes a compromise - the observation could become too unreliable in case of packet reordering and/or severe loss if N is too small, while short flows may not yield a useful upstream loss measurement if N is too large (see [Section 4.2.2](#)).

The value of N should be at least 64 and be a power of 2. This requirement allows an Observer to infer the Q Block length by observing one period of the square signal. It also allows the Observer to identify flows that set the loss bits to arbitrary values (see [Section 7](#)).

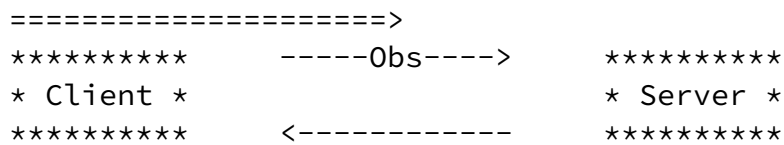
If the sender does not have sufficient information to make an informed decision about Q Block length, the sender should use  $N=64$ , since this value has been extensively tried in large-scale field tests and yielded good results. Alternatively, the sender may also choose a random power-of-2 N for each flow, increasing the chances of using a Q Block length that gives the best signal for some flows.

The sender must keep the value of N constant for a given flow.

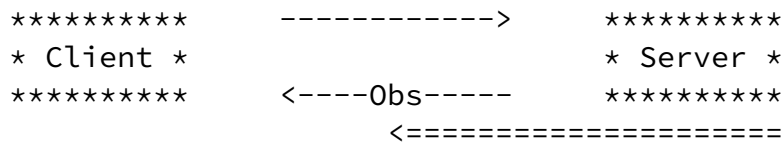
#### 4.2.2. Upstream Loss

Blocks of N (Q Block length) consecutive packets are sent with the same value of the Q bit, followed by another block of N packets with an inverted value of the Q bit. Hence, knowing the value of N, an on-path observer can estimate the amount of upstream loss after observing at least N packets. The upstream loss rate ("uloss") is one minus the average number of packets in a block of packets with the same Q value ("p") divided by N ("uloss=1-avg(p)/N").

The observer needs to be able to tolerate packet reordering that can blur the edges of the square signal, as explained in [Section 4.2.3](#).



(a) in client-server channel (uloss<sub>up</sub>)



(b) in server-client channel (uloss<sub>down</sub>)

Upstream loss

#### 4.2.3. Identifying Q Block Boundaries

Packet reordering can produce spurious edges in the square signal. To address this, the observer should look for packets with the current Q bit value up to X packets past the first packet with a



reverse Q bit value. The value of X, a "Marking Block Threshold", should be less than "N/2".

The choice of X represents a trade-off between resiliency to reordering and resiliency to loss. A very large Marking Block Threshold will be able to reconstruct Q Blocks despite a significant amount of reordering, but it may erroneously coalesce packets from multiple Q Blocks into fewer Q Blocks, if loss exceeds 50% for some Q Blocks.

#### [4.3.](#) L Bit - Loss Event Bit

The Loss Event bit uses an Unreported Loss counter maintained by the protocol that implements the marking mechanism. To use the Loss Event bit, the protocol must allow the sender to identify lost packets. This is true of protocols such as QUIC, partially true for TCP and SCTP (losses of pure ACKs are not detected) and is not true of protocols such as UDP and IP/IPv6.

The Unreported Loss counter is initialized to 0, and L bit of every outgoing packet indicates whether the Unreported Loss counter is positive (L=1 if the counter is positive, and L=0 otherwise).

The value of the Unreported Loss counter is decremented every time a packet with L=1 is sent.

The value of the Unreported Loss counter is incremented for every packet that the protocol declares lost, using whatever loss detection machinery the protocol employs. If the protocol is able to rescind the loss determination later, a positive Unreported Loss counter may be decremented due to the rescission, but it should NOT become negative due to the rescission.

This loss signaling is similar to loss signaling in [[ConEx](#)], except the Loss Event bit is reporting the exact number of lost packets, whereas Echo Loss bit in [[ConEx](#)] is reporting an approximate number of lost bytes.

For protocols, such as TCP ([\[TCP\]](#)), that allow network devices to change data segmentation, it is possible that only a part of the packet is lost. In these cases, the sender must increment Unreported Loss counter by the fraction of the packet data lost (so Unreported

Loss counter may become negative when a packet with L=1 is sent after a partial packet has been lost).

Observation points can estimate the end-to-end loss, as determined by the upstream endpoint, by counting packets in this direction with the L bit equal to 1, as described in [Section 4.3.1](#).

#### [4.3.1](#). End-To-End Loss

The Loss Event bit allows an observer to estimate the end-to-end loss rate by counting packets with L bit value of 0 and 1 for a given flow. The end-to-end loss rate is the fraction of packets with L=1.

The assumption here is that upstream loss affects packets with L=0 and L=1 equally. If some loss is caused by tail-drop in a network device, this may be a simplification. If the sender's congestion controller reduces the packet send rate after loss, there may be a sufficient delay before sending packets with L=1 that they have a greater chance of arriving at the observer.

#### [4.3.2](#). Loss Profile Characterization

In addition to measuring the end-to-end loss rate, the Loss Event bit allows an observer to characterize loss profile, since the distribution of observed packets with L bit set to 1 roughly corresponds to the distribution of packets lost between 1 RTT and 1 RT0 before (see [Section 4.4.1](#)). Hence, observing random single instances of L bit set to 1 indicates random single packet loss, while observing blocks of packets with L bit set to 1 indicates loss affecting entire blocks of packets.

#### [4.4](#). L+Q Bits - Upstream, Downstream, and End-to-End Loss Measurements

Combining L and Q bits allows a passive observer watching a single direction of traffic to accurately measure:

- upstream loss: sender-to-observer loss (see [Section 4.2.2](#))
- downstream loss: observer-to-receiver loss (see [Section 4.4.1.1](#))
- end-to-end loss: sender-to-receiver loss on the observed path (see [Section 4.3.1](#)) with loss profile characterization (see [Section 4.3.2](#))

#### [4.4.1.](#) Correlating End-to-End and Upstream Loss

Upstream loss is calculated by observing packets that did not suffer the upstream loss ([Section 4.2.2](#)). End-to-end loss, however, is calculated by observing subsequent packets after the sender's protocol detected the loss. Hence, end-to-end loss is generally observed with a delay of between 1 RTT (loss declared due to multiple duplicate acknowledgments) and 1 RTO (loss declared due to a timeout) relative to the upstream loss.

The flow RTT can sometimes be estimated by timing protocol handshake messages. This RTT estimate can be greatly improved by observing a dedicated protocol mechanism for conveying RTT information, such as the Spin bit (see [Section 3.1](#)) or Delay bit (see [Section 3.2](#)).

Whenever the observer needs to perform a computation that uses both upstream and end-to-end loss rate measurements, it should use upstream loss rate leading the end-to-end loss rate by approximately 1 RTT. If the observer is unable to estimate RTT of the flow, it should accumulate loss measurements over time periods of at least 4 times the typical RTT for the observed flows.

If the calculated upstream loss rate exceeds the end-to-end loss rate calculated in [Section 4.3.1](#), then either the Q Period is too short for the amount of packet reordering or there is observer loss, described in [Section 4.4.1.2](#). If this happens, the observer should adjust the calculated upstream loss rate to match end-to-end loss rate, unless the following applies.

In case of a protocol like TCP and SCTP that does not track losses of pure ACK packets, observing a direction of traffic dominated by pure ACK packets could result in measured upstream loss that is higher than measured end-to-end loss, if said pure ACK packets are lost upstream. Hence, if the measurement is applied to such protocols, and the observer can confirm that pure ACK packets dominate the observed traffic direction, the observer should adjust the calculated end-to-end loss rate to match upstream loss rate.

##### [4.4.1.1.](#) Downstream Loss

Because downstream loss affects only those packets that did not suffer upstream loss, the end-to-end loss rate ("eloss") relates to

the upstream loss rate ("uloss") and downstream loss rate ("dloss") as  $(1-uloss)(1-dloss)=1-eloss$ . Hence,  $dloss=(eloss-uloss)/(1-uloss)$ .

#### [4.4.1.2](#). Observer Loss

A typical deployment of a passive observation system includes a network tap device that mirrors network packets of interest to a device that performs analysis and measurement on the mirrored packets. The observer loss is the loss that occurs on the mirror path.

Observer loss affects upstream loss rate measurement, since it causes the observer to account for fewer packets in a block of identical Q bit values (see [Section 4.2.2](#)). The end-to-end loss rate measurement, however, is unaffected by the observer loss, since it is a measurement of the fraction of packets with the L bit value of 1, and the observer loss would affect all packets equally (see [Section 4.3.1](#)).

The need to adjust the upstream loss rate down to match end-to-end loss rate as described in [Section 4.4.1](#) is an indication of the observer loss, whose magnitude is between the amount of such adjustment and the entirety of the upstream loss measured in [Section 4.2.2](#). Alternatively, a high apparent upstream loss rate could be an indication of significant packet reordering, possibly due to packets belonging to a single flow being multiplexed over several upstream paths with different latency characteristics.

#### [4.5](#). R Bit - Reflection Square Bit

R bit requires a deployment alongside Q bit. Unlike the square signal for which packets are transmitted into blocks of fixed size, the Reflection square signal (being an alternate marking signal too) produces blocks of packets whose size varies according to these rules:

- when the transmission of a new block starts, its size is set equal to the size of the last Q Block whose reception has been

completed;

- if, before transmission of the block is terminated, the reception of at least one further Q Block is completed, the size of the block is updated to the average size of the further received Q Blocks. Implementation details follow.

The Reflection square value is initialized to 0 and is applied to the R-bit of every outgoing packet. The Reflection square value is toggled for the first time when the completion of a Q Block is detected in the incoming square signal (produced by the opposite node using the Q-bit). When this happens, the number of packets ("p"), detected within this first Q Block, is used to generate a reflection

square signal which toggles every "M=p" packets (at first). This new signal produces blocks of M packets (marked using the R-bit) and each of them is called "Reflection Block" (R Block).

The M value is then updated every time a completed Q Block in the incoming square signal is received, following this formula:  
"M=round(avg(p))".

The parameter "avg(p)" is the average number of packets in a marking period computed considering all the Q Blocks received since the beginning of the current R Block.

To ensure a proper computation of the M value, endpoints implementing the R bit must identify the boundaries of incoming Q Blocks. The same approach described in {#endmarkingblock} should be used.

Looking at the R-bit, unidirectional observation points have an indication of losses experienced by the entire unobserved channel plus those occurred in the path from the sender up to them.

Since the Q Block is sent in one direction, and the corresponding reflected R Block is sent in the opposite direction, the reflected R signal is transmitted with the packet rate of the slowest direction. Namely, if the observed direction is the slowest, there can be multiple Q Blocks transmitted in the unobserved direction before a complete R Block is transmitted in the observed direction. If the unobserved direction is the slowest, the observed direction can be sending R Blocks of the same size repeatedly before it can update the

signal to account for a newly-completed Q Block.

#### [4.5.1.](#) R+Q Bits - Using R and Q Bits for Passive Loss Measurement

Since both sSquare and Reflection square bits are toggled at most every N packets (except for the first transition of the R-bit as explained before), an on-path observer can count the number of packets of each marking block and, knowing the value of N, can estimate the amount of loss experienced by the connection. An observer can calculate different measurements depending on whether it is able to observe a single direction of the traffic or both directions.

Single directional observer:

- upstream loss in the observed direction: the loss between the sender and the observation point (see [Section 4.2.2](#))

- "three-quarters" connection loss: the loss between the receiver and the sender in the unobserved direction plus the loss between the sender and the observation point in the observed direction
- end-to-end loss in the unobserved direction: the loss between the receiver and the sender in the opposite direction

Two directions observer (same metrics seen previously applied to both direction, plus):

- client-observer half round-trip loss: the loss between the client and the observation point in both directions
- observer-server half round-trip loss: the loss between the observation point and the server in both directions
- downstream loss: the loss between the observation point and the receiver (applicable to both directions)

##### [4.5.1.1.](#) Three-Quarters Connection Loss

Except for the very first block in which there is nothing to reflect (a complete Q Block has not been yet received), packets are continuously R-bit marked into alternate blocks of size lower or equal than N. Knowing the value of N, an on-path observer can estimate the amount of loss occurred in the whole opposite channel plus the loss from the sender up to it in the observation channel. As for the previous metric, the "three-quarters" connection loss rate ("tqloss") is one minus the average number of packets in a block of packets with the same R value ("t") divided by "N" ("tqloss=1-avg(t)/N").

```

=====>
= *****      -----Obs----->      *****
= * Client *    <-----              * Server *
= *****      <-----              *****
<=====

```

(a) in client-server channel (tqloss\_up)

```

=====>
*****      ----->      ***** =
* Client *   <-----      * Server * =
*****      <----Obs----- ***** =
<=====

```

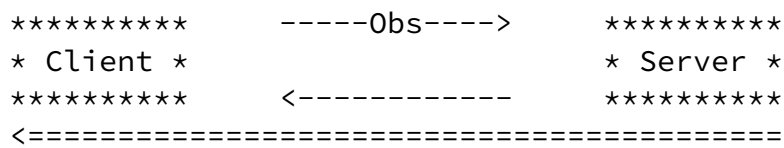
(b) in server-client channel (tqloss\_down)

Three-quarters connection loss

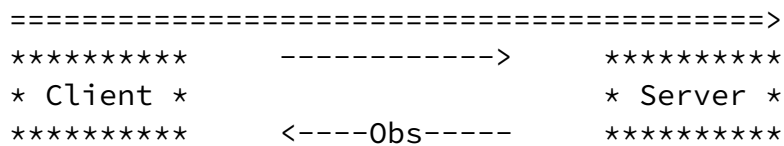
The following metrics derive from this last metric and the upstream loss produced by the Q Bit.

4.5.1.2. End-To-End Loss in the Opposite Direction

End-to-end loss in the unobserved direction ("eloss\_unobserved") relates to the "three-quarters" connection loss ("tqloss") and upstream loss in the observed direction ("uloss") as " $(1-eloss\_unobserved)(1-uloss)=1-tqloss$ ". Hence, " $eloss\_unobserved=(tqloss-uloss)/(1-uloss)$ ".



(a) in client-server channel (eloss\_down)



(b) in server-client channel (eloss\_up)

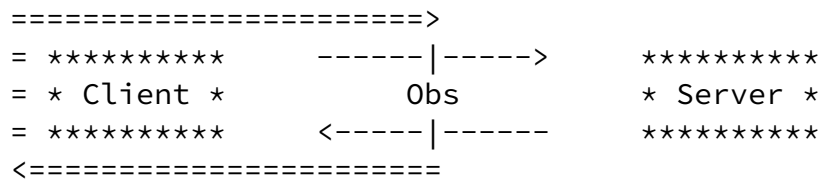
End-To-End loss in the opposite direction

4.5.1.3. Half Round-Trip Loss

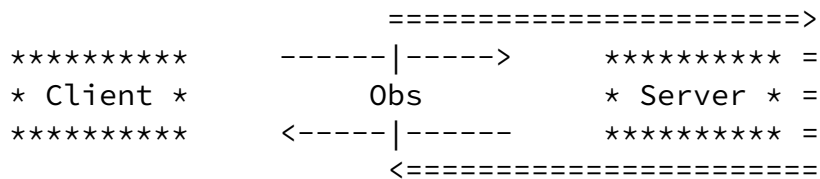
If the observer is able to observe both directions of traffic, it is able to calculate two "half round-trip" loss measurements - loss from the observer to the receiver (in a given direction) and then back to the observer in the opposite direction. For both directions, "half round-trip" loss ("hrtloss") relates to "three-quarters" connection



loss ("tqloss\_opposite") measured in the opposite direction and the upstream loss ("uloss") measured in the given direction as  $(1-uloss)(1-hrtloss)=1-tqloss\_opposite$ . Hence,  $hrtloss=(tqloss\_opposite-uloss)/(1-uloss)$ .



(a) client-observer half round-trip loss (hrtloss\_co)



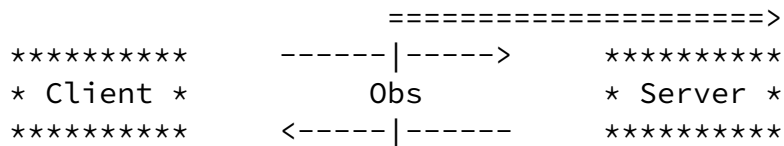
(b) observer-server half round-trip loss (hrtloss\_os)

Half Round-trip loss (both direction)

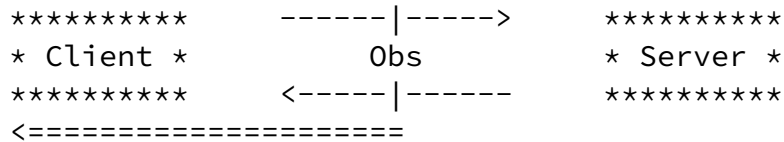
#### [4.5.1.4](#). Downstream Loss

If the observer is able to observe both directions of traffic, it is able to calculate two downstream loss measurements using either end-to-end loss and upstream loss, similar to the calculation in [Section 4.4.1.1](#) or using "half round-trip" loss and upstream loss in the opposite direction.

For the latter,  $dloss=(hrtloss-uloss\_opposite)/(1-uloss\_opposite)$ .



(a) in client-server channel (dloss\_up)



(b) in server-client channel (dloss\_down)

### Downstream loss

#### [4.5.2.](#) Enhancement of R Block Length Computation

The use of the rounding function used in the M computation introduces errors that can be minimized by storing the rounding applied each time M is computed, and using it during the computation of the M value in the following R Block.

This can be achieved introducing the new "r\_avg" parameter in the computation of M. The new formula is "Mr=avg(p)+r\_avg; M=round(Mr); r\_avg=Mr-M" where the initial value of "r\_avg" is equal to 0.

#### [4.5.3.](#) Improved Resilience to Packet Reordering

When a protocol implementing the marking mechanism is able to detect when packets are received out of order, it can improve resilience to packet reordering beyond what is possible using methods described in [Section 4.2.3](#).

This can be achieved by updating the size of the current R Block while this is being transmitted. The reflection block size is then updated every time an incoming reordered packet of the previous Q Block is detected. This can be done if and only if the transmission of the current reflection block is in progress and no packets of the following Q Block have been received.

#### [4.6.](#) Improved Q and R Bits Resilience to Burst Losses

Burst losses can affect Q and R measurements accuracy. Generally, burst losses can be absorbed and correctly measured if smaller than the established Q Block length. On the other hand, entire periods might be wiped out if the burst sizes become too large thus making the observer completely unaware of their loss.

Internet-Draft

Delay and Loss bits

October 2021

To improve burst loss resilience, an observer might consider a received Q or R Block larger than the selected Q Block length as a burst loss event. Then compute the loss as three times Q Block length minus the measured block length. By doing so, an observer can detect burst losses of less than two blocks (e.g., less than 128 packets for Q Block length of 64 packets). A burst loss equal or greater than two consecutive periods would still remain unnoticed by the observer (or underestimated if a period longer than Q Block length were formed).

## 5. Summary of Delay and Loss Marking Methods

This section summarizes the marking methods described in this draft.

For the Delay measurement, it is possible to use the spin bit and/or the delay bit. A unidirectional or bidirectional observer can be used.

Method	# of bits	Available Delay Metrics		Impairments Resiliency	# of meas.
		UNIDIR Observer	BIDIR Observer		
S: Spin Bit	1	RTT	x2 Half RTT	low	very high
D: Delay Bit	1	RTT	x2 Half RTT	high	medium
D <sup>^</sup> : Hidden Delay Bit	1	RTT <sup>^</sup>	x2 Left Half <sup>^</sup> Right Half	high	high
SD: Spin Bit & Delay Bit *	2	RTT	x2 Half RTT	high	very high

x2 Same metric for both directions

\* Both algorithms work independtly; an observer could use approximate spin bit measures when delay bit ones aren't available

<sup>^</sup> Masked metric (real value can be calculated only by those who know

the Additional Delay)

Figure 1: Delay Comparison

For the Loss measurement, each row in the table of Figure 2 represents a loss marking method. For each method the table specifies the number of bits required in the header, the available metrics using an unidirectional or bidirectional observer, applicable protocols, measurement fidelity and delay.

Method	B i t s	Available Loss Metrics		P r o t o c o l s	Measurement Aspects	
		UNIDIR Observer	BIDIR Observer		Fidelity	Delay
T: Round Trip Loss Bit	\$ 1	RT	x2 Half RT	* *	Rate by sampling 1/3 to 1/(3*ppa) of pkts over 2 RTT	~6 RTT
Q: Square Bit	1	Upstream	x2	* *	Rate over N pkts (e.g. 64)	N pkts (e.g. 64)
L: Loss Event Bit	1	E2E	x2	# #	Loss shape (and rate)	Min: RTT Max: RT0
QL: Square + Loss Ev. Bits	2	Upstream Downstream E2E	x2 x2 x2	# # #	-> see Q -> see Q L -> see L	Up: see Q Others: see L
QR: Square + Ref. Sq. Bits	2	Upstream 3/4 RT !E2E	x2 x2 E2E Downstream Half RT	* *	Rate over N*ppa pkts (see Q bit for N)	Up: see Q Others: N*ppa pk (see Q for N)

\* All protocols

# Protocols employing loss detection (w/ or w/o pure ACK loss detection)  
\$ Require a working spin bit  
! Metric relative to the opposite channel  
x2 Same metric for both directions  
ppa Packets-Per-Ack  
Q|L See Q if Upstream loss is significant; L otherwise

Figure 2: Loss Comparison

## [6.](#) ECN-Echo Event Bit

While the primary focus of the draft is on exposing packet loss and delay, modern networks can report congestion before they are forced to drop packets, as described in [\[ECN\]](#). When transport protocols keep ECN-Echo feedback under encryption, this signal cannot be observed by the network operators. When tasked with diagnosing network performance problems, knowledge of a congestion downstream of an observation point can be instrumental.

If downstream congestion information is desired, this information can be signaled with an additional bit.

- E: The "ECN-Echo Event" bit is set to 0 or 1 according to the Unreported ECN Echo counter, as explained below in [Section 6.1](#).

### [6.1.](#) Setting the ECN-Echo Event Bit on Outgoing Packets

The Unreported ECN-Echo counter operates identically to Unreported Loss counter ([Section 4.3](#)), except it counts packets delivered by the network with CE markings, according to the ECN-Echo feedback from the receiver.

This ECN-Echo signaling is similar to ECN signaling in [\[ConEx\]](#). ECN-Echo mechanism in QUIC provides the number of packets received with CE marks. For protocols like TCP, the method described in [\[ConEx-TCP\]](#) can be employed. As stated in [\[ConEx-TCP\]](#), such feedback can be further improved using a method described in [\[ACCURATE\]](#).

## [6.2.](#) Using E Bit for Passive ECN-Reported Congestion Measurement

A network observer can count packets with CE codepoint and determine the upstream CE-marking rate directly.

Observation points can also estimate ECN-reported end-to-end congestion by counting packets in this direction with a E bit equal to 1.

The upstream CE-marking rate and end-to-end ECN-reported congestion can provide information about downstream CE-marking rate. Presence of E bits along with L bits, however, can somewhat confound precise estimates of upstream and downstream CE-markings in case the flow contains packets that are not ECN-capable.

## [7.](#) Protocol Ossification Considerations

Accurate loss and delay information is not critical to the operation of any protocol, though its presence for a sufficient number of flows is important for the operation of networks.

The delay and loss bits are amenable to "greasing" described in [\[RFC8701\]](#), if the protocol designers are not ready to dedicate (and ossify) bits used for loss reporting to this function. The greasing could be accomplished similarly to the Latency Spin bit greasing in [\[QUIC-TRANSPORT\]](#). Namely, implementations could decide that a fraction of flows should not encode loss and delay information and, instead, the bits would be set to arbitrary values. The observers would need to be ready to ignore flows with delay and loss information more resembling noise than the expected signal.

## [8.](#) Examples of Application

### [8.1.](#) QUIC

The binding of a delay signal to QUIC is partially described in [\[QUIC-TRANSPORT\]](#), which adds the spin bit to the first byte of the

short packet header, leaving two reserved bits for future experiments.

To implement the additional signals discussed in this document, the first byte of the short packet header can be modified as follows:

- the delay bit (D) can be placed in the first reserved bit (i.e. the fourth most significant bit `_0x10_`) while the round trip loss bit (T) in the second reserved bit (i.e. the fifth most significant bit `_0x08_`); the proposed scheme is:

```
 0 1 2 3 4 5 6 7
+--+--+--+--+--+--+
|0|1|S|D|T|K|P|P|
+--+--+--+--+--+--+
```

Scheme 1

- alternatively, a two bits loss signal (QL or QR) can be placed in both reserved bits; the proposed schemes, in this case, are:

```
 0 1 2 3 4 5 6 7
+--+--+--+--+--+--+
|0|1|S|Q|L|K|P|P|
+--+--+--+--+--+--+
```

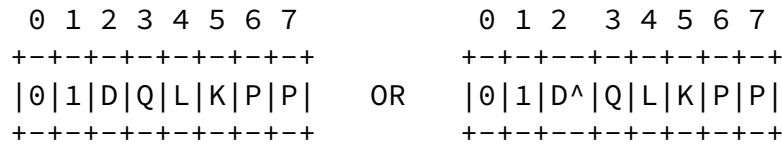
Scheme 2A

```
 0 1 2 3 4 5 6 7
+--+--+--+--+--+--+
|0|1|S|Q|R|K|P|P|
+--+--+--+--+--+--+
```

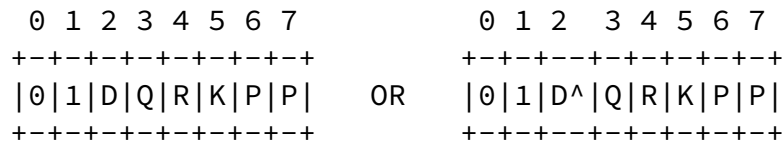
Scheme 2B

A further option would be to substitute the spin bit with the delay

bit (or hidden delay bit) leaving the two reserved bits for loss detection. The proposed schemes are:



Scheme 3A



Scheme 3B

## 8.2. TCP

The signals can be added to TCP by defining bit 4 of byte 13 of the TCP header to carry the spin bit or the delay bit, and possibly bits 5 and 6 to carry additional information, like the delay bit and the round-trip loss bit (DT), or a two bits loss signal (QL or QR).

## 9. Security Considerations

Passive loss and delay observations have been a part of the network operations for a long time, so exposing loss and delay information to the network does not add new security concerns for protocols that are currently observable.

In the absence of packet loss, Q and R bits signals do not provide any information that cannot be observed by simply counting packets

transiting a network path. In the presence of packet loss, Q and R bits will disclose the loss, but this is information about the environment and not the endpoint state. The L bit signal discloses internal state of the protocol's loss detection machinery, but this state can often be gleamed by timing packets and observing congestion controller response.

Hence, loss bits do not provide a viable new mechanism to attack data



integrity and secrecy.

### 9.1. Optimistic ACK Attack

A defense against an Optimistic ACK Attack, described in [\[QUIC-TRANSPORT\]](#), involves a sender randomly skipping packet numbers to detect a receiver acknowledging packet numbers that have never been received. The Q bit signal may inform the attacker which packet numbers were skipped on purpose and which had been actually lost (and are, therefore, safe for the attacker to acknowledge). To use the Q bit for this purpose, the attacker must first receive at least an entire Q Block of packets, which renders the attack ineffective against a delay-sensitive congestion controller.

A protocol that is more susceptible to an Optimistic ACK Attack with the loss signal provided by Q bit and uses a loss-based congestion controller, should shorten the current Q Block by the number of skipped packets numbers. For example, skipping a single packet number will invert the square signal one outgoing packet sooner.

Similar considerations apply to the R Bit, although a shortened R Block along with a matching skip in packet numbers does not necessarily imply a lost packet, since it could be due to a lost packet on the reverse path along with a deliberately skipped packet by the sender.

## 10. Privacy Considerations

To minimize unintentional exposure of information, loss bits provide an explicit loss signal - a preferred way to share information per [\[RFC8558\]](#).

New protocols commonly have specific privacy goals, and loss reporting must ensure that loss information does not compromise those privacy goals. For example, [\[QUIC-TRANSPORT\]](#) allows changing Connection IDs in the middle of a connection to reduce the likelihood of a passive observer linking old and new sub-flows to the same device. A QUIC implementation would need to reset all counters when it changes the destination (IP address or UDP port) or the Connection ID used for outgoing packets. It would also need to avoid

different destination or with a different Connection ID.

## 11. IANA Considerations

This document makes no request of IANA.

## 12. Change Log

TBD

## 13. Contributors

The following people provided valuable contributions to this document:

- Marcus Ihlar, Ericsson, [marcus.ihlar@ericsson.com](mailto:marcus.ihlar@ericsson.com)
- Jari Arkko, Ericsson, [jari.arkko@ericsson.com](mailto:jari.arkko@ericsson.com)
- Emile Stephan, Orange, [emile.stephan@orange.com](mailto:emile.stephan@orange.com)

## 14. Acknowledgements

TBD

## 15. References

### 15.1. Normative References

- [ConEx] Mathis, M. and B. Briscoe, "Congestion Exposure (ConEx) Concepts, Abstract Mechanism, and Requirements", [RFC 7713](#), DOI 10.17487/RFC7713, December 2015, <<https://www.rfc-editor.org/info/rfc7713>>.
- [ConEx-TCP] Kuehlewind, M., Ed. and R. Scheffenegger, "TCP Modifications for Congestion Exposure (ConEx)", [RFC 7786](#), DOI 10.17487/RFC7786, May 2016, <<https://www.rfc-editor.org/info/rfc7786>>.
- [ECN] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", [RFC 3168](#), DOI 10.17487/RFC3168, September 2001, <<https://www.rfc-editor.org/info/rfc3168>>.

- [IP] Postel, J., "Internet Protocol", STD 5, [RFC 791](#), DOI 10.17487/RFC0791, September 1981, <<https://www.rfc-editor.org/info/rfc791>>.
- [IPM-Methods] Morton, A., "Active and Passive Metrics and Methods (with Hybrid Types In-Between)", [RFC 7799](#), DOI 10.17487/RFC7799, May 2016, <<https://www.rfc-editor.org/info/rfc7799>>.
- [IPv6] 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>>.
- [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>>.
- [RFC8558] Hardie, T., Ed., "Transport Protocol Path Signals", [RFC 8558](#), DOI 10.17487/RFC8558, April 2019, <<https://www.rfc-editor.org/info/rfc8558>>.
- [TCP] Postel, J., "Transmission Control Protocol", STD 7, [RFC 793](#), DOI 10.17487/RFC0793, September 1981, <<https://www.rfc-editor.org/info/rfc793>>.

## [15.2](#). Informative References

- [ACCURATE] Briscoe, B., Kuehlewind, M., and R. Scheffenegger, "More Accurate ECN Feedback in TCP", [draft-ietf-tcpm-accurate-ecn-15](#) (work in progress), July 2021.
- [AltMark] Fioccola, G., Ed., Capello, A., Cociglio, M., Castaldelli, L., Chen, M., Zheng, L., Mirsky, G., and T. Mizrahi, "Alternate-Marking Method for Passive and Hybrid Performance Monitoring", [RFC 8321](#), DOI 10.17487/RFC8321, January 2018, <<https://www.rfc-editor.org/info/rfc8321>>.
- [ANRW19-PM-QUIC] Bulgarella, F., Cociglio, M., Fioccola, G., Marchetto, G., and R. Sisto, "Performance measurements of QUIC communications", Proceedings of the Applied Networking Research Workshop, DOI 10.1145/3340301.3341127, July 2019.

Internet-Draft

Delay and Loss bits

October 2021

[I-D.trammell-ippm-spin]

Trammell, B., "An Explicit Transport-Layer Signal for Hybrid RTT Measurement", [draft-trammell-ippm-spin-00](#) (work in progress), January 2019.

[I-D.trammell-tsvwg-spin]

Trammell, B., "A Transport-Independent Explicit Signal for Hybrid RTT Measurement", [draft-trammell-tsvwg-spin-00](#) (work in progress), July 2018.

[IPv6AltMark]

Fioccola, G., Zhou, T., Cociglio, M., Qin, F., and R. Pang, "IPv6 Application of the Alternate Marking Method", [draft-ietf-6man-ipv6-alt-mark-12](#) (work in progress), October 2021.

[QUIC-TRANSPORT]

Iyengar, J., Ed. and M. Thomson, Ed., "QUIC: A UDP-Based Multiplexed and Secure Transport", [RFC 9000](#), DOI 10.17487/RFC9000, May 2021, <<https://www.rfc-editor.org/info/rfc9000>>.

[RFC8517] Dolson, D., Ed., Snellman, J., Boucadair, M., Ed., and C. Jacquenet, "An Inventory of Transport-Centric Functions Provided by Middleboxes: An Operator Perspective", [RFC 8517](#), DOI 10.17487/RFC8517, February 2019, <<https://www.rfc-editor.org/info/rfc8517>>.

[RFC8701] Benjamin, D., "Applying Generate Random Extensions And Sustain Extensibility (GREASE) to TLS Extensibility", [RFC 8701](#), DOI 10.17487/RFC8701, January 2020, <<https://www.rfc-editor.org/info/rfc8701>>.

[RFC9065] Fairhurst, G. and C. Perkins, "Considerations around Transport Header Confidentiality, Network Operations, and the Evolution of Internet Transport Protocols", [RFC 9065](#), DOI 10.17487/RFC9065, July 2021, <<https://www.rfc-editor.org/info/rfc9065>>.

[SPIN-BIT]

Trammell, B., Vaere, P. D., Even, R., Fioccola, G., Fossati, T., Ihlar, M., Morton, A., and E. Stephan, "Adding Explicit Passive Measurability of Two-Way Latency to the QUIC Transport Protocol", [draft-trammell-quic-spin-03](#) (work in progress), May 2018.

Cociglio, et al.

Expires April 28, 2022

[Page 39]

---

Internet-Draft

Delay and Loss bits

October 2021

[UDP-OPTIONS]

Touch, J., "Transport Options for UDP", [draft-ietf-tsvwg-udp-options-13](#) (work in progress), June 2021.

[UDP-SURPLUS]

Herbert, T., "UDP Surplus Header", [draft-herbert-udp-space-hdr-01](#) (work in progress), July 2019.

#### Authors' Addresses

Mauro Cociglio  
Telecom Italia - TIM  
Via Reiss Romoli, 274  
Torino 10148  
Italy

EMail: [mauro.cociglio@telecomitalia.it](mailto:mauro.cociglio@telecomitalia.it)

Alexandre Ferrieux  
Orange Labs

EMail: [alexandre.ferrieux@orange.com](mailto:alexandre.ferrieux@orange.com)

Giuseppe Fioccola  
Huawei Technologies  
Riesstrasse, 25  
Munich 80992  
Germany

EMail: [giuseppe.fioccola@huawei.com](mailto:giuseppe.fioccola@huawei.com)

Igor Lubashev  
Akamai Technologies

EMail: [ilubashe@akamai.com](mailto:ilubashe@akamai.com)

Fabio Bulgarella  
Telecom Italia - TIM  
Via Reiss Romoli, 274  
Torino 10148  
Italy

EMail: [fabio.bulgarella@guest.telecomitalia.it](mailto:fabio.bulgarella@guest.telecomitalia.it)

Cociglio, et al.

Expires April 28, 2022

[Page 40]

---

Internet-Draft

Delay and Loss bits

October 2021

Isabelle Hamchaoui  
Orange Labs

EMail: [isabelle.hamchaoui@orange.com](mailto:isabelle.hamchaoui@orange.com)

Massimo Nilo  
Telecom Italia - TIM  
Via Reiss Romoli, 274  
Torino 10148  
Italy

EMail: [massimo.nilo@telecomitalia.it](mailto:massimo.nilo@telecomitalia.it)

Riccardo Sisto  
Politecnico di Torino

EMail: [riccardo.sisto@polito.it](mailto:riccardo.sisto@polito.it)

Dmitri Tikhonov  
LiteSpeed Technologies

EMail: [dtikhonov@litespeedtech.com](mailto:dtikhonov@litespeedtech.com)

