

Transport Area Working Group
Internet-Draft
Updates: 3819 (if approved)
Intended status: Best Current Practice
Expires: September 4, 2014

B. Briscoe
BT
J. Kaippallimalil
Huawei
P. Thaler
Broadcom Corporation
March 03, 2014

Guidelines for Adding Congestion Notification to Protocols that
Encapsulate IP
draft-briscoe-tsvwg-ecn-encap-guidelines-04

Abstract

The purpose of this document is to guide the design of congestion notification in any lower layer or tunnelling protocol that encapsulates IP. The aim is for explicit congestion signals to propagate consistently from lower layer protocols into IP. Then the IP internetwork layer can act as a portability layer to carry congestion notification from non-IP-aware congested nodes up to the transport layer (L4). Following these guidelines should assure interworking between new lower layer congestion notification mechanisms, whether specified by the IETF or other standards bodies.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on September 4, 2014.

Copyright Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	3
1.1. Scope	5
2. Terminology	6
3. Modes of Operation	7
3.1. Feed-Forward-and-Up Mode	8
3.2. Feed-Up-and-Forward Mode	9
3.3. Feed-Backward Mode	10
3.4. Null Mode	12
4. Feed-Forward-and-Up Mode: Guidelines for Adding Congestion Notification	12
4.1. IP-in-IP Tunnels with Tightly Coupled Shim Headers	13
4.2. Wire Protocol Design: Indication of ECN Support	13
4.3. Encapsulation Guidelines	15
4.4. Decapsulation Guidelines	17
4.5. Sequences of Similar Tunnels or Subnets	18
4.6. Reframing and Congestion Markings	19
5. Feed-Up-and-Forward Mode: Guidelines for Adding Congestion Notification	19
6. Feed-Backward Mode: Guidelines for Adding Congestion Notification	21
7. IANA Considerations (to be removed by RFC Editor)	22
8. Security Considerations	22
9. Conclusions	22
10. Acknowledgements	23
11. Comments Solicited	23
12. References	23
12.1. Normative References	23
12.2. Informative References	24
Appendix A. Outstanding Document Issues	27
Appendix B. Changes in This Version (to be removed by RFC Editor)	27

1. Introduction

The benefits of Explicit Congestion Notification (ECN) described below can only be fully realised if support for ECN is added to the relevant subnetwork technology, as well as to IP. When a lower layer buffer drops a packet obviously it does not just drop at that layer; the packet disappears from all layers. In contrast, when a lower layer marks a packet with ECN, the marking needs to be explicitly propagated up the layers. The same is true if a buffer marks the outer header of a packet that encapsulates inner tunnelled headers. Forwarding ECN is not as straightforward as other headers because it has to be assumed ECN may be only partially deployed. If an egress at any layer is not ECN-aware, or if the ultimate receiver or sender is not ECN-aware, congestion needs to be indicated by dropping a packet, not marking it.

The purpose of this document is to guide the addition of congestion notification to any subnet technology or tunnelling protocol, so that lower layer equipment can signal congestion explicitly and it will propagate consistently into encapsulated (higher layer) headers, otherwise the signals will not reach their ultimate destination.

ECN is defined in the IP header (v4 & v6) [RFC3168] to allow a resource to notify the onset of queue build-up without having to drop packets, by explicitly marking a proportion of packets with the congestion experienced (CE) codepoint.

Given a suitable marking scheme, ECN removes nearly all congestion loss and it cuts delays for two main reasons:

- o It avoids the delay when recovering from congestion losses, which particularly benefits small flows or real-time flows, making their delivery time predictably short [RFC2884];
- o As ECN is used more widely by end-systems, it will gradually remove the need to configure a degree of delay into buffers before they start to notify congestion (the cause of bufferbloat). This is because drop involves a trade-off between sending a timely signal and trying to avoid impairment, whereas ECN is solely a signal not an impairment, so there is no harm triggering it earlier.

Some lower layer technologies (e.g. MPLS, Ethernet) are used to form subnetworks with IP-aware nodes only at the edges. These networks are often sized so that it is rare for interior queues to overflow. However, this has often been more due to the inability of the original TCP protocol to saturate the links. For many years, fixes such as window scaling [RFC1323] proved hard to deploy. But now that modern

operating systems are finally capable of saturating interior links, even the buffers of well-provisioned interior switches will need to signal episodes of queuing.

Propagation of ECN is defined for MPLS [RFC5129], and is being defined for TRILL [trill-rbridge-options], but it remains to be defined for a number of other subnetwork technologies.

Similarly, ECN propagation is yet to be defined for many tunnelling protocols. [RFC6040] defines how ECN should be propagated for IP-in-IP [RFC2003] and IPsec [RFC4301] tunnels. However, as Section 9.3 of RFC3168 pointed out, ECN support will need to be defined for other tunnelling protocols, e.g. L2TP [RFC2661], GRE [RFC1701], [RFC2784], PPTP [RFC2637] and GTP [GTPv1], [GTPv1-U], [GTPv2-C].

Incremental deployment is the most tricky aspect when adding support for ECN. The original ECN protocol in IP [RFC3168] was carefully designed so that a congested buffer would not mark a packet (rather than drop it) unless both source and destination hosts were ECN-capable. Otherwise its congestion markings would never be detected and congestion would just deteriorate further. However, to support congestion marking below the IP layer, it is not sufficient to only check that the two end-points support ECN; correct operation also depends on the decapsulator at each subnet egress faithfully propagating congestion notifications to the higher layer. Otherwise, a legacy decapsulator might silently fail to propagate any ECN signals from the outer to the forwarded header. Then the lost signals would never be detected and again congestion would deteriorate further. The guidelines given later require protocol designers to carefully consider incremental deployment, and suggest various safe approaches for different circumstances.

Of course, the IETF does not have standards authority over every link layer protocol. So this document gives guidelines for designing propagation of congestion notification across the interface between IP and protocols that may encapsulate IP (i.e. that can be layered beneath IP). Each lower layer technology will exhibit different issues and compromises, so the IETF or the relevant standards body must be free to define the specifics of each lower layer congestion notification scheme. Nonetheless, if the guidelines are followed, congestion notification should interwork between different technologies, using IP in its role as a 'portability layer'.

Therefore, the capitalised term 'SHOULD' or 'SHOULD NOT' are often used in preference to 'MUST' or 'MUST NOT', because it is difficult to know the compromises that will be necessary in each protocol design. If a particular protocol design chooses to contradict a

'SHOULD (NOT)' given in the advice below, it MUST include a sound justification.

It has not been possible to give common guidelines for all lower layer technologies, because they do not all fit a common pattern. Instead they have been divided into a few distinct modes of operation: feed-forward-and-upward; feed-upward-and-forward; feed-backward; and null mode. These modes are described in Section 3, then in the following sections separate guidelines are given for each mode.

This document updates the advice to subnetwork designers about ECN in Section 13 of [RFC3819].

1.1. Scope

This document only concerns wire protocol processing of explicit notification of congestion and makes no changes or recommendations concerning algorithms for congestion marking or for congestion response (algorithm issues should be independent of the layer the algorithm operates in).

The question of congestion notification signals with different semantics to those of ECN in IP is touched on in a couple of specific cases (e.g. QCN [IEEE802.1Qau]) and with schemes with multiple severity levels such as PCN [RFC6660]). However, no attempt is made to give guidelines about schemes with different semantics that are yet to be invented.

The semantics of congestion signals can be relative to the traffic class. Therefore correct propagation of congestion signals could depend on correct propagation of any traffic class field between the layers. In this document, correct propagation of traffic class information is assumed, while what 'correct' means and how it is achieved is covered elsewhere (e.g. [RFC2983]) and is outside the scope of the present document.

Note that these guidelines do not require the subnet wire protocol to be changed to accommodate congestion notification. Another way to add congestion notification without consuming header space in the subnet protocol might be to use a parallel control plane protocol.

This document focuses on the congestion notification interface between IP and lower layer protocols that can encapsulate IP, where the term 'IP' includes v4 or v6, unicast, multicast or anycast. However, it is likely that the guidelines will also be useful when a lower layer protocol or tunnel encapsulates itself (e.g. Ethernet MAC in MAC [IEEE802.1Qah]) or when it encapsulates other protocols. In

the feed-backward mode, propagation of congestion signals for multicast and anycast packets is out-of-scope (because it would be so complicated that it is hoped no-one would attempt such an abomination).

2. Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC2119].

Further terminology used within this document:

Protocol data unit (PDU): Information that is delivered as a unit among peer entities of a layered network consisting of protocol control information (typically a header) and possibly user data (payload) of that layer. The scope of this document includes layer 2 and layer 3 networks, where the PDU is respectively termed a frame or a packet (or a cell in ATM). PDU is a general term for any of these. This definition also includes a payload with a shim header lying somewhere between layer 2 & 3.

Transport: The end-to-end transmission control function, conventionally considered at layer-4 in the OSI reference model. Given the audience for this document will often use the word transport to mean low level bit carriage, whenever the term is used it will be qualified, e.g. 'L4 transport'.

Encapsulator: The link or tunnel endpoint function that adds an outer header to a PDU (also termed the 'link ingress', the 'subnet ingress', the 'ingress tunnel endpoint' or just the 'ingress' where the context is clear).

Decapsulator: The link or tunnel endpoint function that removes an outer header from a PDU (also termed the 'link egress', the 'subnet egress', the 'egress tunnel endpoint' or just the 'egress' where the context is clear).

Incoming header: The header of an arriving PDU before encapsulation.

Outer header: The header added to encapsulate a PDU.

Inner header: The header encapsulated by the outer header.

Outgoing header: The header forwarded by the decapsulator.

CE: Congestion Experienced [RFC3168]

ECT: ECN-Capable Transport [RFC3168]

Not-ECT: Not ECN-Capable Transport [RFC3168]

ECN-PDU: A PDU that is part of a feedback loop within which all the nodes that need to propagate explicit congestion notifications back to the Load Regulator are ECN-capable. An IP packet with a non-zero ECN field implies that the endpoints are ECN-capable, so this would be an ECN-PDU. However, ECN-PDU is intended to be a general term for a PDU at any layer, not just IP.

Not-ECN-PDU: A PDU that is part of a feedback-loop within which some nodes necessary to propagate explicit congestion notifications back to the load regulator are not ECN-capable.

Load Regulator: For each flow of PDUs, the transport function that is capable of controlling the data rate. Typically located at the data source, but in-path nodes can regulate load in some congestion control arrangements (e.g. admission control or policing nodes). Note the term "a function capable of controlling the load" deliberately includes a transport that doesn't actually control the load but ideally it ought to (e.g. a sending application without congestion control that uses UDP).

Congestion Baseline: The location of the function on the path that initialised the values of all congestion notification fields in a sequence of packets, before any are set to the congestion experienced (CE) codepoint if they experience congestion further downstream. Typically the original data source at layer-4.

3. Modes of Operation

This section sets down the different modes by which congestion information is passed between the lower layer and the higher one. It acts as a reference framework for the following sections, which give normative guidelines for designers of explicit congestion notification protocols, taking each mode in turn:

Feed-Forward-and-Up: Nodes feed forward congestion notification towards the egress within the lower layer then up and along the layers towards the end-to-end destination at the transport layer. The following local optimisation is possible:

Feed-Up-and-Forward: A lower layer switch feeds-up congestion notification directly into the ECN field in the higher layer (e.g. IP) header, irrespective of whether the node is at the egress of a subnet.

Feed-Backward: Nodes feed back congestion signals towards the ingress of the lower layer and (optionally) attempt to control congestion within their own layer.

Null: Nodes cannot experience congestion at the lower layer except at ingress nodes (which are IP-aware or equivalently higher-layer-aware).

3.1. Feed-Forward-and-Up Mode

Like IP and MPLS, many subnet technologies are based on self-contained protocol data units (PDUs) or frames sent unreliably. They provide no feedback channel at the subnetwork layer, instead relying on higher layers (e.g. TCP) to feed back loss signals.

In these cases, ECN may best be supported by standardising explicit notification of congestion into the lower layer protocol that carries the data forwards. It will then also be necessary to define how the egress of the lower layer subnet propagates this explicit signal into the forwarded upper layer (IP) header. It can then continue forwards until it finally reaches the destination transport (at L4). Then typically the destination will feed this congestion notification back to the source transport using an end-to-end protocol (e.g. TCP). This is the arrangement that has already been used to add ECN to IP-in-IP tunnels [RFC6040], IP-in-MPLS and MPLS-in-MPLS [RFC5129].

This mode is illustrated in Figure 1. Along the middle of the figure, layers 2, 3 & 4 of the protocol stack are shown, and one packet is shown along the bottom as it progresses across the network from source to destination, crossing two subnets connected by a router, and crossing two switches on the path across each subnet. Congestion at the output of the first switch (shown as *) leads to a congestion marking in the L2 header (shown as C in the illustration of the packet). The chevrons show the progress of the resulting congestion indication. It is propagated from link to link across the subnet in the L2 header, then when the router removes the marked L2 header, it propagates the marking up into the L3 (IP) header. The router forwards the marked L3 header into subnet 2, and when it adds a new L2 header it copies the L3 marking into the L2 header as well, as shown by the 'C's in both layers (assuming the technology of subnet 2 also supports explicit congestion marking).

Note that there is no implication that each 'C' marking is encoded the same; a different encoding might be used for the 'C' marking in each protocol.

Finally, for completeness, we show the L3 marking arriving at the destination, where the host transport protocol (e.g. TCP) feeds it

back to the source in the L4 acknowledgement (the 'C' at L4 in the packet at the top of the diagram).

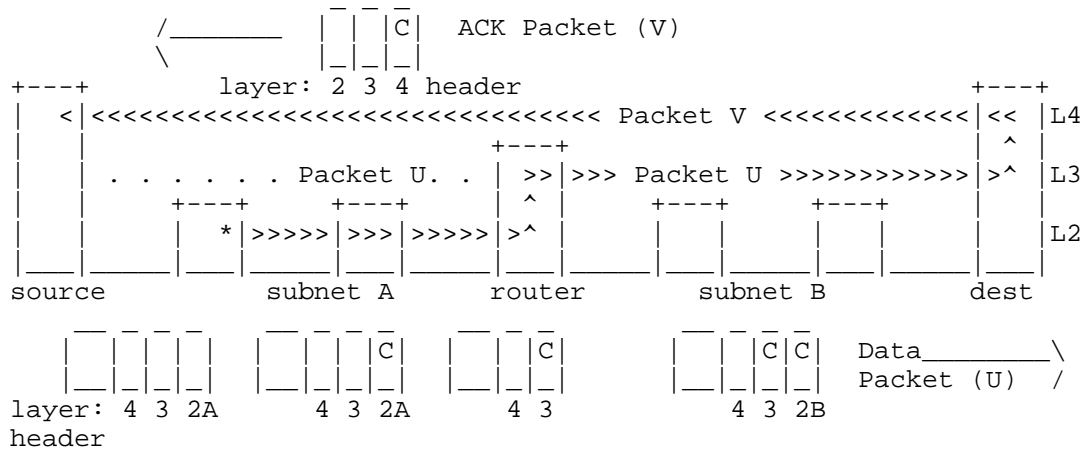


Figure 1: Feed-Forward-and-Up Mode

Of course, modern networks are rarely as simple as this text-book example, often involving multiple nested layers. For example, a 3GPP mobile network may have two IP-in-IP (GTP) tunnels in series and an MPLS backhaul between the base station and the first router. Nonetheless, the example illustrates the general idea of feeding congestion notification forward then upward whenever a header is removed at the egress of a subnet.

Note that the FECN (forward ECN) bit in Frame Relay and the explicit forward congestion indication (EFCI [ITU-T.I.371]) bit in ATM user data cells follow a feed-forward pattern. However, in ATM, this is only as part of a feed-forward-and-backward pattern at the lower layer, not feed-forward-and-up out of the lower layer--the intention was never to interface to IP ECN at the subnet egress. To our knowledge, Frame Relay FECN is solely used to detect where more capacity should be provisioned [Buck00].

3.2. Feed-Up-and-Forward Mode

Ethernet is particularly difficult to extend incrementally to support explicit congestion notification. One way to support ECN in such cases has been to use so called 'layer-3 switches'. These are Ethernet switches that bury into the Ethernet payload to find an IP header and manipulate or act on certain IP fields (specifically Diffserv & ECN). For instance, in Data Center TCP [DCTCP], layer-3 switches are configured to mark the ECN field of the IP header within

the Ethernet payload when their output buffer becomes congested. With respect to switching, a layer-3 switch acts solely on the addresses in the Ethernet header; it doesn't use IP addresses, and it doesn't decrement the TTL field in the IP header.

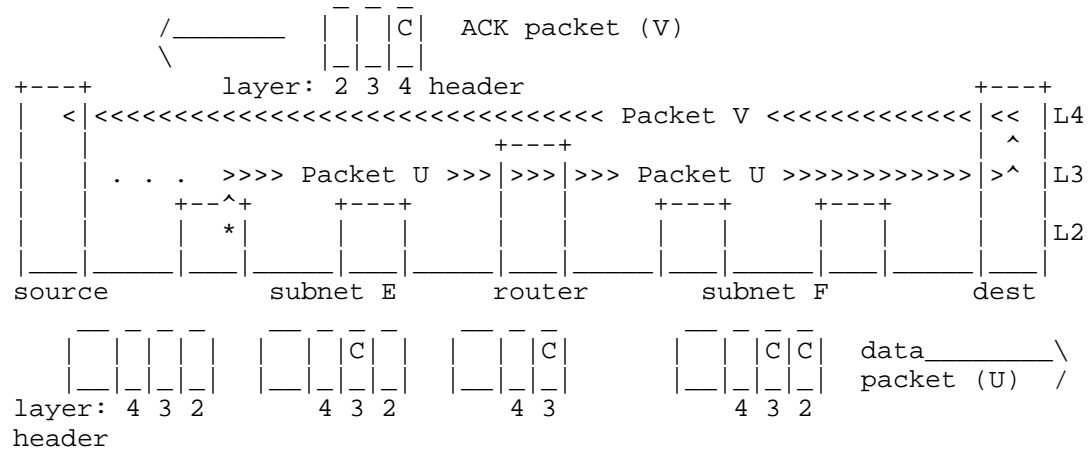


Figure 2: Feed-Up-and-Forward Mode

By comparing Figure 2 with Figure 1, it can be seen that subnet E (perhaps a subnet of layer-3 Ethernet switches) works in feed-up-and-forward mode by notifying congestion directly into L3 at the point of congestion, even though the congested switch does not otherwise act at L3. In this example, the technology in subnet F (e.g. MPLS) does support ECN natively, so when the router adds the layer-2 header it copies the ECN marking from L3 to L2 as well.

3.3. Feed-Backward Mode

In some layer 2 technologies, explicit congestion notification has been defined for use internally within the subnet with its own feedback and load regulation, but typically the interface with IP for ECN has not been defined.

For instance, for the available bit-rate (ABR) service in ATM, the relative rate mechanism was one of the more popular mechanisms for managing traffic, tending to supersede earlier designs. In this approach ATM switches send special resource management (RM) cells in both the forward and backward directions to control the ingress rate of user data into a virtual circuit. If a switch buffer is approaching congestion or congested it sends an RM cell back towards the ingress with respectively the No Increase (NI) or Congestion

Indication (CI) bit set in its message type field [ATM-TM-ABR]. The ingress then holds or decreases its sending bit-rate accordingly.

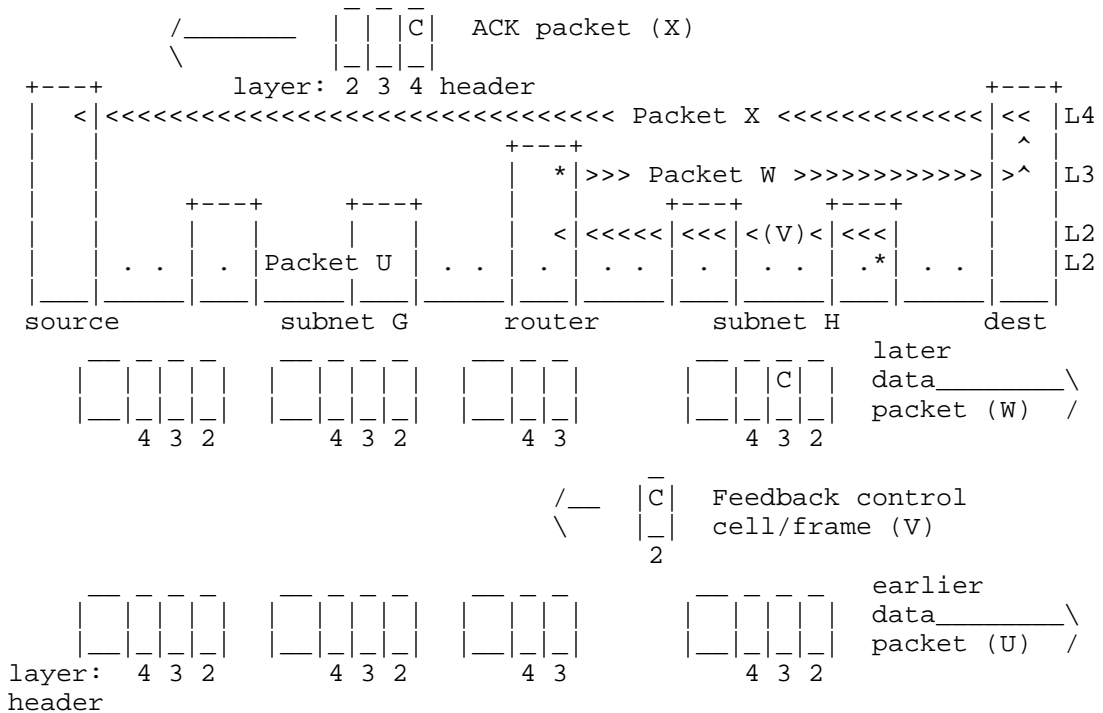


Figure 3: Feed-Backward Mode

ATM's feed-backward approach doesn't fit well when layered beneath IP's feed-forward approach--unless the initial data source is the same node as the ATM ingress. Figure 3 shows the feed-backward approach being used in subnet H. If the final switch on the path is congested (*), it doesn't feed-forward any congestion indications on packet (U). Instead it sends a control cell (V) back to the router at the ATM ingress.

However, the backward feedback doesn't reach the original data source directly because IP doesn't support backward feedback (and subnet G is independent of subnet H). Instead, the router in the middle throttles down its sending rate but the original data sources don't reduce their rates. The resulting rate mismatch causes the middle router's buffer at layer 3 to back up until it becomes congested, which it signals forwards on later data packets at layer 3 (e.g. packet W). Note that the forward signal from the middle router is not triggered directly by the backward signal. Rather, it is

triggered by congestion resulting from the middle router's mismatched rate response to the backward signal.

In response to this later forward signalling, end-to-end feedback at layer-4 finally completes the tortuous path of congestion indications back to the origin data source, as before.

3.4. Null Mode

Often link and physical layer resources are 'non-blocking' by design. In these cases congestion notification may be implemented but it does not need to be deployed at the lower layer; ECN in IP would be sufficient.

A degenerate example is a point-to-point Ethernet link. Excess loading of the link merely causes the queue from the higher layer to back up, while the lower layer remains immune to congestion. Even a whole meshed subnetwork can be made immune to interior congestion by limiting ingress capacity and careful sizing of links, particularly if multi-path routing is used to ensure even worst-case patterns of load cannot congest any link.

4. Feed-Forward-and-Up Mode: Guidelines for Adding Congestion Notification

Feed-forward-and-up is the mode already used for signalling ECN up the layers through MPLS into IP [RFC5129] and through IP-in-IP tunnels [RFC6040]. These RFCs take a consistent approach and the following guidelines are designed to ensure this consistency continues as ECN support is added to other protocols that encapsulate IP. The guidelines are also designed to ensure compliance with the more general best current practice for the design of alternate ECN schemes given in [RFC4774].

The rest of this section is structured as follows:

- o Section 4.1 addresses the most straightforward cases, where [RFC6040] can be applied directly to add ECN to tunnels that are effectively the same as IP-in-IP tunnels.
- o The subsequent sections give guidelines for adding ECN to a subnet technology that uses feed-forward-and-up mode like IP, but it is not so similar to IP that [RFC6040] rules can be applied directly. Specifically:
 - * Sections 4.2, 4.3 and 4.4 respectively address how to add ECN support to the wire protocol and to the encapsulators and decapsulators at the ingress and egress of the subnet.

- * Section 4.5 deals with the special, but common, case of sequences of tunnels or subnets that all use the same technology
- * Section 4.6 deals with the question of reframing when IP packets do not map 1:1 into lower layer frames.

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

A common pattern for many tunnelling protocols is to encapsulate an inner IP header with shim header(s) then an outer IP header. In many cases the shim header(s) always have to be tightly coupled to the outer IP header because they are not sufficient as outer headers in their own right. In such cases the shim header(s) and the outer IP header are always added (or removed) in the same operation. Therefore, in all such tightly coupled IP-in-IP tunnelling protocols, the rules in [RFC6040] for propagating the ECN field between the two IP headers SHOULD be applied directly.

Examples of tightly coupled IP-in-IP tunnelling protocols where [RFC6040] can be applied directly are:

- o L2TP [RFC2661]
- o GRE [RFC1701], [RFC2784]
- o PPTP [RFC2637]
- o GTP [GTPv1], [GTPv1-U], [GTPv2-C]
- o VXLAN [vxlan].

4.2. Wire Protocol Design: Indication of ECN Support

This section is intended to guide the redesign of any lower layer protocol that encapsulate IP to add native ECN support at the lower layer. It reflects the approaches used in [RFC6040] and in [RFC5129]. Therefore IP-in-IP tunnels or IP-in-MPLS or MPLS-in-MPLS encapsulations that already comply with [RFC6040] or [RFC5129] will already satisfy this guidance.

A lower layer (or subnet) congestion notification system:

1. SHOULD NOT apply explicit congestion notifications to PDUs that are destined for legacy layer-4 transport implementations that will not understand ECN, and

2. SHOULD NOT apply explicit congestion notifications to PDUs if the egress of the subnet might not propagate congestion notifications onward into the higher layer.

We use the term ECN-PDUs for a PDU on a feedback loop that will propagate congestion notification properly because it meets both the above criteria. And a Not-ECN-PDU is a PDU on a feedback loop that does not meet both criteria, and will therefore not propagate congestion notification properly. A corollary of the above is that a lower layer congestion notification protocol:

3. SHOULD be able to distinguish ECN-PDUs from Not-ECN-PDUs.

Note that there is no need for all interior nodes within a subnet to be able to mark congestion explicitly. A mix of ECN and drop signals from different nodes is fine. However, if any interior nodes might generate ECN markings, guideline 2 above says that all relevant egress node(s) SHOULD be able to propagate those markings up to the higher layer.

In IP, if the ECN field in each PDU is cleared to the Not-ECT (not ECN-capable transport) codepoint, it indicates that the L4 transport will not understand congestion markings. A congested buffer must not mark these Not-ECT PDUs, and therefore drops them instead.

The mechanism a lower layer uses to distinguish the ECN-capability of PDUs need not mimic that of IP. All the above guidelines say is that the lower layer system, as a whole, should achieve the same outcome. For instance, ECN-capable feedback loops might use PDUs that are identified by a particular set of labels or tags. Alternatively, logical link protocols that use flow state might determine whether a PDU can be congestion marked by checking for ECN-support in the flow state. Other protocols might depend on out-of-band control signals.

The per-domain checking of ECN support in MPLS [RFC5129] is a good example of a way to avoid sending congestion markings to transports that will not understand them, without using any header space in the subnet protocol.

In MPLS, header space is extremely limited, therefore RFC5129 does not provide a field in the MPLS header to indicate whether the PDU is an ECN-PDU or a Not-ECN-PDU. Instead, interior nodes in a domain are allowed to set explicit congestion indications without checking whether the PDU is destined for a transport that will understand them. Nonetheless, this is made safe by requiring that the network operator upgrades all decapsulating edges of a whole domain at once, as soon as even one switch within the domain is configured to mark rather than drop during congestion. Therefore, any edge node that

might decapsulate a packet will be capable of checking whether the higher layer transport is ECN-capable. When decapsulating a CE-marked packet, if the decapsulator discovers that the higher layer (inner header) indicates the transport is not ECN-capable, it drops the packet--effectively on behalf of the earlier congested node (see Decapsulation Guideline 1 in Section 4.4).

It was only appropriate to define such an incremental deployment strategy because MPLS is targeted solely at professional operators, who can be expected to ensure that a whole subnetwork is consistently configured. This strategy might not be appropriate for other link technologies targeted at zero-configuration deployment or deployment by the general public (e.g. Ethernet). For such 'plug-and-play' environments it will be necessary to invent a failsafe approach that ensures congestion markings will never fall into black holes, no matter how inconsistently a system is put together. Alternatively, congestion notification relying on correct system configuration could be confined to flavours of Ethernet intended only for professional network operators, such as IEEE 802.1ah Provider Backbone Bridges (PBB).

QCN [IEEE802.1Qau] provides another example of how to indicate to lower layer devices that the end-points will not understand ECN. An operator can define certain 802.1p classes of service to indicate non-QCN frames and an ingress bridge is required to map arriving not-QCN-capable IP packets to one of these non-QCN 802.1p classes.

4.3. Encapsulation Guidelines

This section is intended to guide the redesign of any node that encapsulates IP with a lower layer header when adding native ECN support to the lower layer protocol. It reflects the approaches used in [RFC6040] and in [RFC5129]. Therefore IP-in-IP tunnels or IP-in-MPLS or MPLS-in-MPLS encapsulations that already comply with [RFC6040] or [RFC5129] will already satisfy this guidance.

1. Egress Capability Check: A subnet ingress needs to be sure that the corresponding egress of a subnet will propagate any congestion notification added to the outer header across the subnet. This is necessary in addition to checking that an incoming PDU indicates an ECN-capable (L4) transport. Examples of how this guarantee might be provided include:
 - * by configuration (e.g. if any label switches in a domain support ECN marking, [RFC5129] requires all egress nodes to have been configured to propagate ECN)

- * by the ingress explicitly checking that the egress propagates ECN (e.g. TRILL uses IS-IS to check path capabilities before using critical options [trill-rbridge-options])
 - * by inherent design of the protocol (e.g. by encoding ECN marking on the outer header in such a way that a legacy egress that does not understand ECN will consider the PDU corrupt and discard it, thus at least propagating a form of congestion signal).
2. Egress Fails Capability Check: If the ingress cannot guarantee that the egress will propagate congestion notification, the ingress SHOULD disable ECN when it forwards the PDU at the lower layer. An example of how the ingress might disable ECN at the lower layer would be by setting the outer header of the PDU to identify it as a Not-ECN-PDU, assuming the subnet technology supports such a concept.
 3. Standard Congestion Monitoring Baseline: Once the ingress to a subnet has established that the egress will correctly propagate ECN, on encapsulation it SHOULD encode the same level of congestion in outer headers as is arriving in incoming headers. For example it might copy any incoming congestion notification into the outer header of the lower layer protocol.

This ensures that all outer headers reflect congestion accumulated along the whole upstream path since the Load Regulator, not just since the ingress of the subnet. A node that is not the Load Regulator SHOULD NOT re-initialise the level of CE markings in the outer to zero.

This guideline is intended to ensure that any bulk congestion monitoring of outer headers (e.g. by a network management node monitoring ECN in passing frames) is most meaningful. For instance, if an operator measures CE in 0.4% of passing outer headers, this information is only useful if the operator knows where the proportion of CE markings was last initialised to 0% (the Congestion Baseline). Such monitoring information will not be useful if some subnet ingress nodes reset all outer CE markings while others copy incoming CE markings into the outer.

Most information can be extracted if the Congestion Baseline is standardised at the node that is regulating the load (the Load Regulator--typically the data source). Then the operator can measure both congestion since the Load Regulator, and congestion since the subnet ingress. The latter might be measurable by subtracting the level of CE markings on inner headers from that on outer headers (see Appendix C of [RFC6040]).

4.4. Decapsulation Guidelines

This section is intended to guide the redesign of any node that decapsulates IP from within a lower layer header when adding native ECN support to the lower layer protocol. It reflects the approaches used in [RFC6040] and in [RFC5129]. Therefore IP-in-IP tunnels or IP-in-MPLS or MPLS-in-MPLS encapsulations that already comply with [RFC6040] or [RFC5129] will already satisfy this guidance.

A subnet egress SHOULD NOT simply copy congestion notification from outer headers to the forwarded header. It SHOULD calculate the outgoing congestion notification field from the inner and outer headers using the following guidelines. If there is any conflict, rules earlier in the list take precedence over rules later in the list:

1. If the arriving inner header is a Not-ECN-PDU it implies the L4 transport will not understand explicit congestion markings.
Then:
 - * If the outer header carries an explicit congestion marking, the packet SHOULD be dropped--the only indication of congestion that the L4 transport will understand.
 - * If the outer is an ECN-PDU that carries no indication of congestion or a Not-ECN-PDU the PDU SHOULD be forwarded, but still as a Not-ECN-PDU.
2. If the outer header does not support explicit congestion notification (a Not-ECN-PDU), but the inner header does (an ECN-PDU), the inner header SHOULD be forwarded unchanged.
3. In some lower layer protocols congestion may be signalled as a numerical level, such as in the control frames of quantised congestion notification [IEEE802.1Qau]. If such a multi-bit encoding encapsulates an ECN-capable IP data packet, a function will be needed to convert the quantised congestion level into the frequency of congestion markings in outgoing IP packets.
4. Congestion indications may be encoded by a severity level. For instance increasing levels of congestion might be encoded by numerically increasing indications, e.g. pre-congestion notification (PCN) can be encoded in each PDU at three severity levels in IP or MPLS [RFC6660].

If the arriving inner header is an ECN-PDU, where the inner and outer headers carry indications of congestion of different

severity, the more severe indication SHOULD be forwarded in preference to the less severe.

5. The inner and outer headers might carry a combination of congestion notification fields that should not be possible given any currently used protocol transitions. For instance, if Encapsulation Guideline 3 in Section 4.3 had been followed, it should not be possible to have a less severe indication of congestion in the outer than in the inner. It MAY be appropriate to log unexpected combinations of headers and possibly raise an alarm.

If a safe outgoing codepoint can be defined for such a PDU, the PDU SHOULD be forwarded rather than dropped. Some implementers discard PDUs with currently unused combinations of headers just in case they represent an attack. However, an approach using alarms and policy-mediated drop is preferable to hard-coded drop, so that operators can keep track of possible attacks but currently unused combinations are not precluded from future use through new standards actions.

4.5. Sequences of Similar Tunnels or Subnets

In some deployments, particularly in 3GPP networks, an IP packet may traverse two or more IP-in-IP tunnels in sequence that all use identical technology (e.g. GTP).

In such cases, it would be sufficient for every encapsulation and decapsulation in the chain to comply with RFC 6040. Alternatively, as an optimisation, a node that decapsulates a packet and immediately re-encapsulates it for the next tunnel MAY copy the incoming outer ECN field directly to the outgoing outer and the incoming inner ECN field directly to the outgoing inner. Then the overall behavior across the sequence of tunnel segments would still be consistent with RFC 6040.

Appendix C of RFC6040 describes how a tunnel egress can monitor how much congestion has been introduced within a tunnel. A network operator might want to monitor how much congestion had been introduced within a whole sequence of tunnels. Using the technique in Appendix C of RFC6040 at the final egress, the operator could monitor the whole sequence of tunnels, but only if the above optimisation were used consistently along the sequence of tunnels, in order to make it appear as a single tunnel. Therefore, tunnel endpoint implementations SHOULD allow the operator to configure whether this optimisation is enabled.

When ECN support is added to a subnet technology, consideration SHOULD be given to a similar optimisation between subnets in sequence if they all use the same technology.

4.6. Reframing and Congestion Markings

The guidance in this section is worded in terms of framing boundaries, but it applies equally whether the protocol data units are frames, cells or packets.

Where framing boundaries are different between two layers, congestion indications SHOULD be propagated on the basis that a congestion indication on a PDU applies to all the octets in the PDU. On average, an encapsulator or decapsulator SHOULD approximately preserve the number of marked octets arriving and leaving (counting the size of inner headers, but not added encapsulating headers).

The next departing frame SHOULD be immediately marked even if only enough incoming marked octets have arrived for part of the departing frame. This ensures that any outstanding congestion marked octets are propagated immediately, rather than held back waiting for a frame no bigger than the outstanding marked octets--which might involve a long wait.

For instance, an algorithm for marking departing frames could maintain a counter representing the balance of arriving marked octets minus departing marked octets. It adds the size of every marked frame that arrives and if the counter is positive it marks the next frame to depart and subtracts its size from the counter. This will often leave a negative remainder in the counter, which is deliberate.

5. Feed-Up-and-Forward Mode: Guidelines for Adding Congestion Notification

The guidance in this section is applicable when IP packets:

- o are encapsulated in Ethernet headers;
- o are forwarded by the eNode-B (base station) of a 3GPP radio access network, which is required to apply ECN marking during congestion [LTE-RA].

This guidance also generalises to encapsulation by other subnet technologies with no native support for explicit congestion notification at the lower layer, but with support for finding and processing an IP header. It is unlikely to be applicable or necessary for IP-in-IP encapsulation, where feed-forward-and-up mode based on [RFC6040] would be more appropriate.

Marking the IP header while switching at layer-2 (by using a layer-3 switch) or while forwarding in a radio access network seems to represent a layering violation. However, it can be considered as a benign optimisation if the guidelines below are followed. Feed-up-and-forward is certainly not a general alternative to implementing feed-forward congestion notification in the lower layer, because:

- o IPv4 and IPv6 are not the only layer-3 protocols that might be encapsulated by lower layer protocols
- o Link-layer encryption might be in use, making the layer-2 payload inaccessible
- o Many Ethernet switches do not have 'layer-3 switch' capabilities so they cannot read or modify an IP payload
- o It might be costly to find an IP header (v4 or v6) when it may be encapsulated by more than one lower layer header, e.g. Ethernet MAC in MAC [IEEE802.1Qah].

Nonetheless, configuring lower layer equipment to look for an ECN field in an encapsulated IP header is a useful optimisation. If the implementation follows the guidelines below, this optimisation does not have to be confined to a controlled environment such as within a data centre; it could usefully be applied on any network--even if the operator is not sure whether the above issues will never apply:

1. If a native lower-layer congestion notification mechanism exists for a subnet technology, it is safe to mix feed-up-and-forward with feed-forward-and-up on other switches in the same subnet. However, it will generally be more efficient to use the native mechanism.
2. The depth of the search for an IP header SHOULD be limited. If an IP header is not found soon enough, or an unrecognised or unreadable header is encountered, the switch SHOULD resort to an alternative means of signalling congestion (e.g. drop, or the native lower layer mechanism if available).
3. It is sufficient to use the first IP header found in the stack; the egress of the relevant tunnel can propagate congestion notification upwards to any more deeply encapsulated IP headers later.

6. Feed-Backward Mode: Guidelines for Adding Congestion Notification

It can be seen from Section 3.3 that congestion notification in a subnet using feed-backward mode has generally not been designed to be directly coupled with IP layer congestion notification. The subnet attempts to minimise congestion internally, and if the incoming load at the ingress exceeds the capacity somewhere through the subnet, the layer 3 buffer into the ingress backs up. Thus, a feed-backward mode subnet is in some sense similar to a null mode subnet, in that there is no need for any direct interaction between the subnet and higher layer congestion notification. Therefore no detailed protocol design guidelines are appropriate. Nonetheless, a more general guideline is appropriate:

1. A subnetwork technology intended to eventually interface to IP SHOULD NOT be designed using only the feed-backward mode, which is certainly best for a stand-alone subnet, but would need to be modified to work efficiently as part of the wider Internet, because IP uses feed-forward-and-up mode.

The feed-backward approach at least works beneath IP, where the term 'works' is used only in a narrow functional sense because feed-backward can result in very inefficient and sluggish congestion control--except if it is confined to the subnet directly connected to the original data source, when it is faster than feed-forward. It would be valid to design a protocol that could work in feed-backward mode for paths that only cross one subnet, and in feed-forward-and-up mode for paths that cross subnets.

In the early days of TCP/IP, a similar feed-backward approach was tried for explicit congestion signalling, using source-quench (SQ) ICMP control packets. However, SQ fell out of favour and is now formally deprecated [RFC6633]. The main problem was that it is hard for a data source to tell the difference between a spoofed SQ message and a quench request from a genuine buffer on the path. It is also hard for a lower layer buffer to address an SQ message to the original source port number, which may be buried within many layers of headers, and possibly encrypted.

Quantised congestion notification (QCN--also known as backward congestion notification or BCN) [IEEE802.1Qau] uses a feed-backward mode structurally similar to ATM's relative rate mechanism. However, QCN confines its applicability to scenarios such as some data centres where all endpoints are directly attached by the same Ethernet technology. If a QCN subnet were later connected into a wider IP-based internetwork (e.g. when attempting to interconnect multiple data centres) it would suffer the inefficiency shown Figure 3.

7. IANA Considerations (to be removed by RFC Editor)

This memo includes no request to IANA.

8. Security Considerations

If a lower layer wire protocol is redesigned to include explicit congestion signalling in-band in the protocol header, care SHOULD be taken to ensure that the field used is specified as mutable during transit. Otherwise interior nodes signalling congestion would invalidate any authentication protocol applied to the lower layer header--by altering a header field that had been assumed as immutable.

The redesign of protocols that encapsulate IP in order to propagate congestion signals between layers raises potential signal integrity concerns. Experimental or proposed approaches exist for assuring the end-to-end integrity of in-band congestion signals, e.g.:

- o Congestion exposure (ConEx) for networks to audit that their congestion signals are not being suppressed by other networks or by receivers, and for networks to police that senders are responding sufficiently to the signals, irrespective of the transport protocol used [I-D.ietf-conex-abstract-mech].
- o The ECN nonce [RFC3540] for a TCP sender to detect whether a network or the receiver is suppressing congestion signals.
- o A test with the same goals as the ECN nonce, but without the need for the receiver to co-operate with the protocol [I-D.moncaster-tcpm-rcv-cheat].

Given these end-to-end approaches are already being specified, it would make little sense to attempt to design hop-by-hop congestion signal integrity into a new lower layer protocol, because end-to-end integrity inherently achieves hop-by-hop integrity.

9. Conclusions

Following the guidance in the document enables ECN support to be extended to numerous protocols that encapsulate IP (v4 & v6) in a consistent way, so that IP continues to fulfil its role as an end-to-end interoperability layer. This includes:

- o A wide range of tunnelling protocols with various forms of shim header between two IP headers;

- o A wide range of subnet technologies, particularly those that work in the same 'feed-forward-and-up' mode that is used to support ECN in IP and MPLS.

Guidelines have been defined for supporting propagation of ECN between Ethernet and IP on so-called Layer-3 Ethernet switches, using a 'feed-up-and-forward' mode. This approach could enable other subnet technologies to pass ECN signals into the IP layer, even if they do not support ECN natively.

Finally, attempting to add ECN to a subnet technology in feed-backward mode is deprecated except in special cases, due to its likely sluggish response to congestion.

10. Acknowledgements

Thanks to Gorry Fairhurst for extensive reviews. Thanks also to the following reviewers: Ingemar Johansson and Piers O'Hanlon and Michael Welzl, who pointed out that lower layer congestion notification signals may have different semantics to those in IP.

Bob Briscoe was part-funded by the European Community under its Seventh Framework Programme through the Trilogy project (ICT-216372) for initial drafts and through the Reducing Internet Transport Latency (RITE) project (ICT-317700) subsequently. The views expressed here are solely those of the authors.

11. Comments Solicited

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

12. References

12.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.
- [RFC3819] Karn, P., Bormann, C., Fairhurst, G., Grossman, D., Ludwig, R., Mahdavi, J., Montenegro, G., Touch, J., and L. Wood, "Advice for Internet Subnetwork Designers", BCP 89, RFC 3819, July 2004.

- [RFC4774] Floyd, S., "Specifying Alternate Semantics for the Explicit Congestion Notification (ECN) Field", BCP 124, RFC 4774, November 2006.
- [RFC5129] Davie, B., Briscoe, B., and J. Tay, "Explicit Congestion Marking in MPLS", RFC 5129, January 2008.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, November 2010.

12.2. Informative References

- [ATM-TM-ABR] Cisco, "Understanding the Available Bit Rate (ABR) Service Category for ATM VCs", Design Technote 10415, June 2005.
- [Buck00] Buckwalter, J., "Frame Relay: Technology and Practice", Pub. Addison Wesley ISBN-13: 978-0201485240, 2000.
- [DCTCP] Alizadeh, M., Greenberg, A., Maltz, D., Padhye, J., Patel, P., Prabhakar, B., Sengupta, S., and M. Sridharan, "Data Center TCP (DCTCP)", ACM SIGCOMM CCR 40(4)63--74, October 2010, <<http://portal.acm.org/citation.cfm?id=1851192>>.
- [GTPv1-U] 3GPP, "General Packet Radio System (GPRS) Tunnelling Protocol User Plane (GTPv1-U)", Technical Specification TS 29.281, .
- [GTPv1] 3GPP, "GPRS Tunnelling Protocol (GTP) across the Gn and Gp interface", Technical Specification TS 29.060, .
- [GTPv2-C] 3GPP, "Evolved General Packet Radio Service (GPRS) Tunnelling Protocol for Control plane (GTPv2-C)", Technical Specification TS 29.274, .
- [I-D.ietf-conex-abstract-mech] Mathis, M. and B. Briscoe, "Congestion Exposure (ConEx) Concepts and Abstract Mechanism", draft-ietf-conex-abstract-mech-08 (work in progress), October 2013.
- [I-D.moncaster-tcpm-rcv-cheat] Moncaster, T., "A TCP Test to Allow Senders to Identify Receiver Non-Compliance", draft-moncaster-tcpm-rcv-cheat-01 (work in progress), June 2007.

[IEEE802.1Qah]

IEEE, "IEEE Standard for Local and Metropolitan Area Networks--Virtual Bridged Local Area Networks--Amendment 6: Provider Backbone Bridges", IEEE Std 802.1Qah-2008, August 2008, <<http://www.ieee802.org/1/pages/802.1ah.html>>.

(Access Controlled link within page)

[IEEE802.1Qau]

Finn, N., Ed., "IEEE Standard for Local and Metropolitan Area Networks--Virtual Bridged Local Area Networks - Amendment 13: Congestion Notification", IEEE Std 802.1Qau-2010, March 2010, <<http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=5454061>>.

(Access Controlled link within page)

[ITU-T.I.371]

ITU-T, "Traffic Control and Congestion Control in B-ISDN", ITU-T Rec. I.371 (03/04), March 2004, <<http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=5454061>>.

[LTE-RA] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2", Technical Specification TS 36.300, .

[RFC1323] Jacobson, V., Braden, B., and D. Borman, "TCP Extensions for High Performance", RFC 1323, May 1992.

[RFC1701] Hanks, S., Li, T., Farinacci, D., and P. Traina, "Generic Routing Encapsulation (GRE)", RFC 1701, October 1994.

[RFC2003] Perkins, C., "IP Encapsulation within IP", RFC 2003, October 1996.

[RFC2637] Hamzeh, K., Pall, G., Verthein, W., Taarud, J., Little, W., and G. Zorn, "Point-to-Point Tunneling Protocol", RFC 2637, July 1999.

[RFC2661] Townsley, W., Valencia, A., Rubens, A., Pall, G., Zorn, G., and B. Palter, "Layer Two Tunneling Protocol "L2TP"", RFC 2661, August 1999.

- [RFC2784] Farinacci, D., Li, T., Hanks, S., Meyer, D., and P. Traina, "Generic Routing Encapsulation (GRE)", RFC 2784, March 2000.
- [RFC2884] Hadi Salim, J. and U. Ahmed, "Performance Evaluation of Explicit Congestion Notification (ECN) in IP Networks", RFC 2884, July 2000.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", RFC 2983, October 2000.
- [RFC3540] Spring, N., Wetherall, D., and D. Ely, "Robust Explicit Congestion Notification (ECN) Signaling with Nonces", RFC 3540, June 2003.
- [RFC4301] Kent, S. and K. Seo, "Security Architecture for the Internet Protocol", RFC 4301, December 2005.
- [RFC6633] Gont, F., "Deprecation of ICMP Source Quench Messages", RFC 6633, May 2012.
- [RFC6660] Briscoe, B., Moncaster, T., and M. Menth, "Encoding Three Pre-Congestion Notification (PCN) States in the IP Header Using a Single Diffserv Codepoint (DSCP)", RFC 6660, July 2012.
- [trill-rbridge-options] Eastlake, D., Ghanwani, A., Manral, V., and C. Bestler, "RBridges: Further TRILL Header Extensions", draft-ietf-trill-rbridge-options-07 (work in progress), June 2012.
- [vxlan] Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger, L., Sridhar, T., Bursell, M., and C. Wright, "VXLAN: A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks", draft-mahalingam-dutt-dcops-vxlan-08 (work in progress), February 2014.

Appendix A. Outstanding Document Issues

1. [GF] Concern that certain guidelines warrant a MUST (NOT) rather than a SHOULD (NOT). Given the guidelines say that if any SHOULD (NOT)s are not followed, a strong justification will be needed, they have been left as SHOULD (NOT) pending further list discussion. In particular:
 - * If inner is a Not-ECN-PDU and Outer is CE (or highest severity congestion level), MUST (not SHOULD) drop?
2. Consider whether an IETF Standard Track doc will be needed to Update the IP-in-IP protocols listed in Section 4.1--at least those that the IET

Appendix B. Changes in This Version (to be removed by RFC Editor)

From briscoe-03 to 04:

- * Re-arranged the introduction to describe the purpose of the document first before introducing ECN in more depth. And clarified the introduction throughout.
- * Added applicability to 3GPP TS 36.300.

From briscoe-02 to 03:

- * Scope section:
 - + Added dependence on correct propagation of traffic class information
 - + For the feed-backward mode, deemed multicast and anycast out of scope
- * Ensured all guidelines referring to subnet technologies also refer to tunnels and vice versa by adding applicability sentences at the start of sections 4.1, 4.2, 4.3, 4.4, 4.6 and 5.
- * Added Security Considerations on ensuring congestion signal fields are classed as immutable and on using end-to-end congestion signal integrity technologies rather than hop-by-hop.

From briscoe-01 to 02:

- * Added authors: JK & PT

- * Added
 - + Section 4.1 "IP-in-IP Tunnels with Tightly Coupled Shim Headers"
 - + Section 4.5 "Sequences of Similar Tunnels or Subnets"
 - + roadmap at the start of Section 4, given the subsections have become quite fragmented.
 - + Section 9 "Conclusions"
- * Clarified why transports are starting to be able to saturate interior links
- * Under Section 1.1, addressed the question of alternative signal semantics and included multicast & anycast.
- * Under Section 3.1, included a 3GPP example.
- * Section 4.2. "Wire Protocol Design":
 - + Altered guideline 2. to make it clear that it only applies to the immediate subnet egress, not later ones
 - + Added a reminder that it is only necessary to check that ECN propagates at the egress, not whether interior nodes mark ECN
 - + Added example of how QCN uses 802.1p to indicate support for QCN.
- * Added references to Appendix C of RFC6040, about monitoring the amount of congestion signals introduced within a tunnel
- * Appendix A: Added more issues to be addressed, including plan to produce a standards track update to IP-in-IP tunnel protocols.
- * Updated acks and references

From briscoe-00 to 01:

- * Intended status: BCP (was Informational) & updates 3819 added.
- * Briefer Introduction: Introductory para justifying benefits of ECN. Moved all but a brief enumeration of modes of operation

to their own new section (from both Intro & Scope). Introduced incr. deployment as most tricky part.

- * Tightened & added to terminology section
- * Structured with Modes of Operation, then Guidelines section for each mode.
- * Tightened up guideline text to remove vagueness / passive voice / ambiguity and highlight main guidelines as numbered items.
- * Added Outstanding Document Issues Appendix
- * Updated references

Authors' Addresses

Bob Briscoe
BT
B54/77, Adastral Park
Martlesham Heath
Ipswich IP5 3RE
UK

Phone: +44 1473 645196
EMail: bob.briscoe@bt.com
URI: <http://bobbbriscoe.net/>

John Kaippallimalil
Huawei
5340 Legacy Drive, Suite 175
Plano, Texas 75024
USA

EMail: john.kaippallimalil@huawei.com

Pat Thaler
Broadcom Corporation
5025 Keane Drive
Carmichael, CA 95608
USA

EMail: pthaler@broadcom.com

Internet Engineering Task Force
Internet-Draft
Intended status: Informational
Expires: April 24, 2014

T. Eckert, Ed.
R. Penno
A. Choukir
C. Eckel
Cisco Systems, Inc.
October 21, 2013

A Framework for Signaling Flow Characteristics between Applications and
the Network
draft-eckert-intarea-flow-metadata-framework-02

Abstract

This document provides a framework for communicating information elements (a.k.a. metadata) in a consistent manner between applications and the network to provide better visibility of application flows, thereby enabling differentiated treatment of those flows. These information elements can be conveyed using various signaling protocols, including PCP, RSVP, and STUN.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on April 24, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect

to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	3
2. Background	3
2.1. Deep packet inspection	4
2.1.1. Benefits	4
2.1.2. Limitation	4
2.2. Explicit signaling methods	5
3. Proposed framework	6
3.1. Overview	7
3.1.1. Common, application independent, IPFIX registered, information elements	7
3.1.2. Cross-protocol information element encoding rules . .	7
3.1.3. Anticipated Usage Models	8
3.1.3.1. Informational	8
3.1.3.2. Advisory	8
3.1.3.3. Service Request	9
3.1.4. Considerations for signaling of common information elements	9
3.1.4.1. Proxy originated information	9
3.1.4.2. Authentication	9
3.1.4.3. Common encoding	10
3.1.4.4. Usage Model to Protocol integration	10
3.2. Proposed common information elements	11
3.2.1. Bandwidth Attributes	12
3.2.1.1. Maximum Bandwidth	12
3.2.1.2. Minimum Bandwidth	12
3.2.1.3. Bandwidth Pool	12
3.2.2. Traffic Class Attributes	12
3.2.2.1. RFC4594-DSCP	12
3.2.2.2. Traffic Class Label (TCL)	12
3.2.3. Acceptable Path Attributes	13
3.2.3.1. Delay Tolerance	13
3.2.3.2. Loss Tolerance	14
3.2.3.3. Jitter Tolerance	14
3.2.4. Application Identification	15
3.2.4.1. RFC 6759 style application identification	15
3.2.4.2. URL style application identification	15
4. Acknowledgements	17
5. Informative References	17
Authors' Addresses	18

1. Introduction

This document provides a framework for communicating information elements (a.k.a. metadata) in a consistent manner between applications and the network to provide better visibility of application flows, thereby enabling differentiated treatment of those flows. These information elements can be conveyed using various signaling protocols, including PCP, RSVP, and STUN.

The framework is built around the definition of four key components:

1. A set of application independent information elements (IEs)
2. An encoding of these IEs that is independent of the signaling protocol used as transport
3. Usages of these IEs to support various transactional semantics
4. A mapping of one or more of these usages to an initial set of signaling protocols, including PCP, RSVP, and STUN

This document defines an initial set of IEs, a set of encoding rules, and initial usage model. The actual encoding is defined in [I-D.choukir-tsvwg-flow-metadata-encoding]. Additional documents define the mapping to specific signaling protocols (e.g. RSVP [I-D.zamfir-tsvwg-flow-metadata-rsvp], STUN [I-D.martinsen-mmusic-malice], and PCP [I-D.wing-pcp-flowdata])

2. Background

This section provides background on the motivation for the framework.

Identification and treatment of application flows are critical for the successful deployment and operation of applications based on a wide range of signaling protocols. Historically, this functionality has been accomplished to the extent possible using heuristics, which inspect and infer flow characteristics.

Heuristics may be based on port ranges, IP subnetting, or deep packet inspection (DPI), e.g. application level gateway (ALG). Port based solutions suffer from port overloading and inconsistent port usage. IP subnetting solutions are error prone and result in network management hassle. DPI is computationally expensive and becomes a challenge with the wider adoption of encrypted signaling and secured traffic. An additional drawback of DPI is that the resulting insights are not available, or need to be recomputed, at network nodes further down the application flow path.

The proposed solution allows applications to explicitly signal their flow characteristics to the network. It also provides network nodes with visibility of the application flow characteristics and enables them to contribute to the flow description. The resulting flow description may be communicated as feedback from the network to applications.

The proposed solution does not enhance existing heuristic based mechanisms, nor does it preclude the use of such mechanisms. Rather, it proposes a new mechanism that does not suffer the drawbacks of heuristic based mechanisms.

2.1. Deep packet inspection

2.1.1. Benefits

Deep Packet Inspection (DPI) and other traffic observation methods (such as performance monitoring) are successfully being used for two type of workflows:

1. Provide network operators with visibility into traffic for troubleshooting, capacity planning, accounting and billing and other off network workflows. This is done by exporting observed traffic analysis via protocol such as IPFIX and SNMP.
2. Provide differentiated network services for the traffic according to network operator defined rule sets, including policing and shaping of traffic, providing admission control, impacting routing, permitting passage of traffic (e.g. firewall functions), etc.

Note: For the context of this document, we consider that DPI starts as early into packets as using ACLs with UDP/TCP port numbers to classify traffic.

2.1.2. Limitation

These two workflows, visibility and differentiated network services, are critical in many networks. However, their reliance on inspection and observation limits the ability to enable these workflows more widely.

- o Simple observation based classification, especially ones relying on TCP/UDP, ports often result in incorrect results due to port overloading (i.e. ports used by applications other than those claiming the port with IANA).

- o More and more traffic is encrypted, rendering deep packet inspection impossible or much more complex (e.g. needing to share encryption keys with network equipment).
- o Observation generally requires inspecting the control and signaling traffic of applications. This traffic may flow through a different network path than the actual application data traffic. Impacting the traffic behavior is ineffective in those scenarios.
- o Observation of control, signaling and data traffic with DPI will in general result in less insight into the applications intent than if the application was explicitly signaling its intent to the network.
- o Without explicit desire by the application to signal its intent to the network, it will also not consider to explicitly provide authentication to the network. DPI mechanism have a more difficult job in analyzing application traffic when authentication mechanisms are in use (if they even can)
- o Without explicit involvement of the application, network services leveraging DPI traffic classification impact the application behavior by impacting its traffic, but cannot provide explicit feedback to the application in the form of signaling.

2.2. Explicit signaling methods

There are a variety of existing and evolving signaling options that can provide explicit application to network signaling and serve the visibility and differentiated network services workflows where DPI is currently being used. It seems clear that there will be no single one-protocol-fits-all solution. Every protocol is currently defined in its own silo, creating duplicate or inconsistent information models. This results in duplicate work, more operational complexity and an inability to easily convert information between protocols to easily leverage the best protocol option for each specific use case. Examples of existing signaling options include the following:

- o RSVP is the original on path signaling protocol standardized by the IETF. It operates on path out-of-band and could support any transport protocol traffic (it currently supports TCP and UDP). Its original goal was to provide admission control. Arguably, its success was impacted by its reliance on router-alert because this often leads to RSVP packets being filtered by intervening networks. To date, more lightweight signaling workflows utilizing RSVP have not been standardized within the IETF.

- o NSIS (next Steps in Signaling) is the next iteration of RSVP-like signaling defined by the IETF. Because it focused on the same fundamental workflow as RSVP admission control as its main driver, and because it did not provide significant enough use-case benefits over RSVP, it has seen even less adoption than RSVP.
- o STUN is an on path, in-band signaling protocol that could easily be extended to provide signaling to on path network devices because it provides an easily inspected packet signature, at least for transport protocols such as UDP and SCTP. Through its extensions TURN and ICE, it is becoming quite popular in application signaling driven by the initial use-case of automatically opening up firewall pinholes and determining the best local and remote addresses for peer-to-peer connectivity (ICE).
- o PCP is a protocol designed to support use cases similar to UPnP firewall traversal. It also can easily be extended to provide more generic application to network signaling for traffic flows. Unlike the prior protocols, it is not meant to be used on path end-to-end but rather independently on one "edge" of a traffic flow. It is therefore an attractive alternative (albeit with challenges under path redundancy) because it allows the introduction of application to network signaling without relying on the remote peer. This is especially useful in multi-domain communications.
- o In addition to these, depending on the devices where it is performed, different degrees of DPI may be used to achieve explicit signaling. For example, inspection of HTTP connections is often viable in high-touch network devices. Such inspection may provide explicit signaling if the application purposely keeps or inserts information elements that are meant to be signaled to the network in the clear, or knowingly uses an encryption scheme shared with the network.

Rather than encourage independent, protocol specific solutions to this problem, this document provides a protocol and application independent framework that can be applied in a consistent fashion across the various protocols.

3. Proposed framework

3.1. Overview

The proposed framework includes the following elements:

3.1.1. Common, application independent, IPFIX registered, information elements

An application media flow may be expressed as a set of information elements that are defined and registered like observation-based IPFIX attributes. We propose leveraging IPFIX as the information model (not necessarily as the transport signaling) for the following reasons:

- o As outlined above, export of traffic information is one of the two big workflows. IPFIX is arguably the most flexible, extensible and best defined option for this. Leveraging the same information model for flow characteristics facilitates export of this information via IPFIX.
- o IPFIX allows for IETF/IANA standardized information elements, but also for unambiguous vendor-defined attributes by including the so-called PEN (Private Enterprise Number) into the information element type. Note that IPFIX has ongoing work to better disseminate vendor specific registration of attributes. The framework defined here expects to be able to leverage the output of that work.

3.1.2. Cross-protocol information element encoding rules

The majority of the protocols listed previously (RSVP, NSIS, STUN/ICE, PCP) require (or favor) compact binary encoding of information elements. This is natively supported by the information element registration of IPFIX.

The IPFIX registry defines each information element's data-type, and there is a native binary network encoding for each of these types. At a minimum, every protocol leveraging common information elements would need to use an encoding that identifies the information element's PEN and IE-ID, and that leverages network standard binary encoding of the value including the length of the value. Including the length of the value into the encoding is required for extensibility because otherwise new information elements could not be introduced without first having all network devices know the data-type, and therefore the length, of the information element. Leveraging network standard binary encoding is equally important to permit network elements to propagate information elements from one protocol to another protocol without understanding the information elements data-type.

In protocols that are not constrained to binary encoding, it is nevertheless highly desirable to include the equivalent information and therefore permit propagation between binary and non-binary transport of information elements without having to understand all information elements.

3.1.3. Anticipated Usage Models

The signaling of information elements may be from application to the network or from network to application. When signaled within a given protocol, the information elements may be interpreted independently of that protocol, or it may be used in combination with the given protocol.

3.1.3.1. Informational

The most simplistic usage model is one in which applications signal information elements describing their anticipated or existing flows into the network along the path of those flows without expecting or requiring anything back from the network. Network elements along the flow path may or may not do something with this information.

This "informational" usage model enables network elements along the path to support the workflows traditionally performed via DPI mechanisms, as described previously.

3.1.3.2. Advisory

This usage model extends the "informational" usage in that the application expects or requests some information back from the network. With this usage, the same information elements apply and may be communicated by the application into the network, but the application indicates its interest in receiving some feedback.

Default values are defined for each information element to unambiguously support cases in which an application does not have a valid value to communicate with the network; rather, it wants the network to provide a value back to it in response. In essence, this allows an application to ask a question and receive an answer from the network. Of course, a network element may provide similar feedback for cases in which an application communicated a non-default value as well. Network elements may also provide unsolicited advisory feedback.

In all cases, applications are not guaranteed to receive an answer or any specific service from the network. In the event an answer is provided, that answer is similarly not a guarantee of any specific service or treatment by the network. It is to be interpreted as advisory only.

As mentioned previously, the same information elements are used in the signaling from the application to the network as well as from the network to the application. The underlying transport protocol used to carry the information elements is expected to provide the necessary request/response semantics or some other mechanism by which the communication in both directions can be tied together.

3.1.3.3. Service Request

This usage model extends the "advisory" usage to operate as an explicit service request. Unlike the advisory usage, information elements signalled by the application are interpreted by network elements within the context of a service request, and information elements signalled by the network back to the application are interpreted within the context of a response to that request.

As with the advisory usage, the same information elements are used in the signaling from the application to the network as well as from the network to the application. The underlying transport protocol used to carry the information elements is expected to provide the necessary service request/response semantics.

3.1.4. Considerations for signaling of common information elements

3.1.4.1. Proxy originated information

The goal of this framework is to enable applications to explicitly signal common information elements about their traffic flows and optionally receive common information elements from the network as feedback. Nevertheless, it is clear that broad adoption of such technology is improved by enabling the use of proxies. The proxies can provide or amend the flow description information in the absence of Flow Metadata support by the application itself.

3.1.4.2. Authentication

Common information elements should provide for cryptographic authentication by the sender. In general the authentication provides some form of identification of the sender and proves that the common information elements covered by the authentication were originated from, or approved by, that identity.

3.1.4.3. Common encoding

A companion document [I-D.choukir-tsvwg-flow-metadata-encoding] covers recommended encoding rules that take the following aspects into account:

- o Compact binary encoding rules
- o Signaling for both sent and received traffic flows
- o Signaling of standard and vendor specific information elements
- o Minimizes protocol specific definition required to add informational or advisory common information elements into existing transactions
- o Signaling of feedback from the network
- o Identification of originator to support proxies and facilitate mitigation between common information elements from different originators
- o Signaling of authenticators

3.1.4.4. Usage Model to Protocol integration

There is a range of options for how this framework is integrated with a particular transport protocol. We describe two examples we consider useful:

3.1.4.4.1. Common transport informative integration

1. A transport protocol signaling method is defined to carry the common encoded information elements at least in signaling from application to network.
2. If the transport by itself does not already have a mechanism to indicate a purely informative protocol transaction, then a protocol specific indication for this is added.

In result, this integration achieves two option:

1. Informative common information elements can be sent from application to network by using the protocol's method to indicate the purely informational protocol transaction. This option effectively leverages the protocol as transport for additional informative attribute based services without impacting the services and transactions of the protocol otherwise.

2. Informative common information elements can be sent alongside an existing protocol transaction. In this case they may either be ships in the night (triggering informative attribute based services), or they may additionally be used by the policy rules of the protocol transaction itself which could be advisory or service request. All feedback of the transaction would still rely on protocol specific information element (common information elements only used from host to network).

This integration is for example defined in [I-D.wing-pcp-flowdata], [I-D.zamfir-tsvwg-flow-metadata-rsvp], and [I-D.martinsen-mmusic-malice].

3.1.4.4.2. Common transport advisory integration

In addition to the common transport informative integration, the transport encoding is extended to carry the common transport information element in feedback messages from the network to the host /application. The method to indicate informative only transaction, when sending to the network is used to indicate advisory only transaction when signaling from the network.

This option primarily enables informative and advisory usage models, but it can equally interact with pre-existing service-request options of the transport protocol and impact advisory feedback or the service request itself based on that interaction.

3.2. Proposed common information elements

The section defines an initial set of common information elements. These information elements are intended to be added to the set of IANA standardized information elements either by this or associated documents. Additional documents are expected to define additional attributes that can use either IANA or other vendor-PEN.

All information element definition must include the following:

1. Default value to be provided by an application when it does not have an informative value to provide to the network, but is interested in receiving an advisory value of the attribute from the network. If no advisory feedback is requested, and no informative value is known, the attribute may simply not be sent.
2. Conflict resolution in the presence of different values for the same information element (e.g. two peers signaling information elements for both the upstream and downstream direction of a flow include different values for the information element)

3.2.1. Bandwidth Attributes

3.2.1.1. Maximum Bandwidth

This attribute is used to convey the maximum sustained bandwidth for the flow. It is an unsigned 64 bit value and is specified in bits per second.

Default Value: 0

Conflict Resolution: Minimum for the set of non-default values

3.2.1.2. Minimum Bandwidth

This attribute is used to convey the minimum sustained bandwidth for the flow. It is an unsigned 64 bit value and is specified in bits per second. Not sending the Minimum Bandwidth is equivalent to sending the same value as for Maximum Bandwidth.

Default Value: 0

Conflict Resolution: Minimum of the set of non-default values

3.2.1.3. Bandwidth Pool

This attribute is used to convey that the traffic dynamically shares bandwidth with other traffic using the same Bandwidth Pool. Variable length GUID (Global Unique ID) of at least 48 bits. The Maximum Bandwidth used by the pool is the largest Maximum Bandwidth indicated by any member, the Minimum Bandwidth of the Pool is the largest Minimum Bandwidth indicated by any member.

3.2.2. Traffic Class Attributes

3.2.2.1. RFC4594-DSCP

This attribute is used to convey the DSCP value appropriate for the flow. It is an unsigned 8 bit value. Values signaled are assumed to be in compliance with [RFC4594] or backward compatible extensions thereof. Other values are undefined.

Default Value: 0xff

Conflict Resolution: tbd

3.2.2.2. Traffic Class Label (TCL)

The data type of this information element is a string. It carries the Traffic Class Label defined in [I-D.ietf-mmusic-traffic-class-for-sdp]. Depending on the outcome of that drafts standardization, the version carried as an information element may be slightly expanded over the its definition for SDP. The TCL is a structured string of the form:

`<category>.<application>(.adjective)(.adjective)`

category and application provide a base categorization of the traffic class that attempts to provide a simplified and extensible, framework for the traffic class definitions in [RFC4594]. These base classifications can be refined with zero or more adjectives. Examples of a TCL is "conversational.video.avconf".

Default Value: Empty string

Conflict Resolution: tbd

3.2.3. Acceptable Path Attributes

The following set of attributes deal with tolerance to various path impairments. A discrete and ordered set of values is defined for each. This way the values are applicable on a per hop basis as well as end to end. The values may be mapped to relevant metrics within a given network, such as the mapping of delay tolerance and loss tolerance to QCI values as defined in [I-D.penno-pcp-mobile-qos]

3.2.3.1. Delay Tolerance

This attribute is used to convey the delay tolerance of an application with respect to the associated flow. When provided by a network element, it indicates the delay tolerance expected of the application with respect to the associated flow. It is a 3 bit field for which values are assigned as follows:

0 = no information available

1 = very low

2 = low

3 = medium

4 = high

5-7: reserved

Default Value: 0

Conflict Resolution: For application to network, the most strict of non-default values. For network to application, the least strict of the set of non-default values.

3.2.3.2. Loss Tolerance

This attribute is used to convey the loss tolerance of an application with respect to the associated flow. When provided by a network element, it indicates the loss tolerance expected of the application with respect to the associated flow. It is a 3 bit field for which values are assigned as follows:

0 = no information available

1 = very low

2 = low

3 = medium

4 = high

5-7: reserved

Default Value: 0

Conflict Resolution: For application to network, the most strict of non-default values. For network to application, the least strict of the set of non-default values.

3.2.3.3. Jitter Tolerance

This attribute is used to convey the jitter tolerance of an application with respect to the associated flow. When provided by a network element, it indicates the jitter tolerance expected of the application with respect to the associated flow. It is a 3 bit field for which values are assigned as follows:

0 = no information available

1 = very low

2 = low

3 = medium

4 = high

5-7: reserved

Default Value: 0

Conflict Resolution: For application to network, the most strict of non-default values. For network to application, the least strict of the set of non-default values.

3.2.4. Application Identification

Application identification is clearly one of the more difficult classification goals. The proposals included here are as of yet not widely vetted:

3.2.4.1. RFC 6759 style application identification

[RFC6759] defines the IPFIX IE-IDs that permit both IANA and vendor specific application identification. Though defined for observation (a.k.a.: DPI), it could also be used with explicit signaling from applications.

Applications that use one of the protocols for which there is an IANA port allocation could explicitly indicate this port via the IANA-L4 engine-id in their application to network signaling. This would identify the application even if the application is not using the IANA assigned port for it. This covers cases in which applications use ports other than registered, such as HTTP servers running on other than 80, or when ports get mapped due to PAT.

To avoid collision with DPI exported IANA-L4 classification, it is necessary to assign a new engine-id for application-self assigned IANA-L4 classification (e.g. new engine-id for IANA-L4-SELF-ASSIGNED). If an application vendor has a PEN, the application can use a PANA-L7-PEN classification with the PEN of the originating application vendor. Likewise, if applications are in general made available via "market" type reseller mechanism (common in mobile device applications), then the application vendor could request an application identification from the market owner and leverage the market owners PEN.

3.2.4.2. URL style application identification

One problem with [RFC6759] style application identification especially non-IANA registered ones is the complexity in making all network elements learn the semantic of the numeric encoding of e.g. the PANA-L7-PEN information element in signaling protocols that only

use the numeric encoding of information elements. The second problem may be to determine what PEN to use, because not every developer of an application may be a company that has a PEN or otherwise would intend to apply for one. Application identification via a URL encoded string information element is a way to overcome both issues. Today, almost all applications have some DNS domain associated with them through which they are being marketed or that belongs to the company developing the application. Therefore, one simple form of self assigned application identification is a new IPFIX information element: `UrlAppId`. The value of this information element is an abbreviated URL of the following form:

```
<fqdn> / <app-name> /[ <version> | <other-details> ]
```

The idea is that the owner of `<fqdn>` (fully qualified domain name) is assigning an `<app-name>`, and by signaling both `<domain-name>` and `<app-name>`, this information element provides a self-identifying, unambiguous application identification.

Example:

```
example.com/network-lemmings/sdn-edition
```

A game publishing house or application market operator with the domain name `example.com` is initially allocating the `UrlAppId` `example.com/network-lemmings` to that application. After 35 years, a new variant of the game is released, the SDN edition, and the app-developer decides that it would best like to distinguish this application variant by the above `UrlAppId` `example.com/network-lemmings/sdn-edition`.

In general, different traffic flows within a single application should best not be distinguished via the `UrlAppId`, but instead rely on attributes more specifically targeted for that purpose (such as the `TrafficClassLabel`). If there is no adequate better attribute defined, application developers may choose to use the other-details section of the `UrlAppId` to distinguish flows within the same application.

Formally, the only requirement against the `UrlAppId` is that the `fqdn` part is a DNS domain owned by the assigner, and that the rest of the string after the first `/` is as self explanatory as possible.

It should be noted that in the context of DPI, classification of web-based application traffic is very often performed by URL inspection of HTTP traffic. This proposed intent based information element leverages that model and makes it usable where it can not be currently used with just DPI: encrypted HTTP, non-HTTP applications, HTTP applications with non-descriptive URLs, etc.

4. Acknowledgements

The authors would like to thank Dan Wing, Anca Zamfir, Paul Jones, and Tirumaleswar Reddy for their valuable contributions to this document.

5. Informative References

- [I-D.choukir-tsvwg-flow-metadata-encoding]
Eckert, T., Zamfir, A., Choukir, A., and C. Eckel,
"Protocol Independent Encoding for Signaling Flow
Characteristics", draft-choukir-tsvwg-flow-metadata-
encoding-01 (work in progress), July 2013.
- [I-D.ietf-mmusic-traffic-class-for-sdp]
Polk, J., Dhesikan, S., and P. Jones, "The Session
Description Protocol (SDP) 'trafficclass' Attribute",
draft-ietf-mmusic-traffic-class-for-sdp-04 (work in
progress), July 2013.
- [I-D.martinsen-mmusic-malice]
Penno, R., Martinsen, P., Wing, D., and A. Zamfir, "Meta-
data Attribute signaling with ICE", draft-martinsen-
mmusic-malice-00 (work in progress), July 2013.
- [I-D.penno-pcp-mobile-qos]
Penno, R., Reddy, T., Wing, D., Steeg, B., and M.
Boucadair, "PCP Usage for Quality of Service (QoS) in
Mobile Networks", draft-penno-pcp-mobile-qos-00 (work in
progress), July 2013.
- [I-D.wing-pcp-flowdata]
Wing, D., Penno, R., and T. Reddy, "PCP Flowdata Option",
draft-wing-pcp-flowdata-00 (work in progress), July 2013.
- [I-D.zamfir-tsvwg-flow-metadata-rsvp]
Eckert, T., Zamfir, A., and A. Choukir, "Flow Metadata
Signaling with RSVP", draft-zamfir-tsvwg-flow-metadata-
rsvp-00 (work in progress), July 2013.

[RFC4594] Babiarz, J., Chan, K., and F. Baker, "Configuration Guidelines for DiffServ Service Classes", RFC 4594, August 2006.

[RFC6759] Claise, B., Aitken, P., and N. Ben-Dvora, "Cisco Systems Export of Application Information in IP Flow Information Export (IPFIX)", RFC 6759, November 2012.

Authors' Addresses

Toerless Eckert (editor)
Cisco Systems, Inc.
San Jose
US

Email: eckert@cisco.com

Reinaldo Penno
Cisco Systems, Inc.
170 West Tasman Drive
San Jose 95134
USA

Email: repenno@cisco.com

Amine Choukir
Cisco Systems, Inc.
Lausanne
CH

Email: amchouki@cisco.com

Charles Eckel
Cisco Systems, Inc.
170 West Tasman Drive
San Jose, CA 95134
US

Email: eckelcu@cisco.com

TSVWG
Internet-Draft
Intended status: Informational
Expires: May 16, 2015

R. Geib, Ed.
Deutsche Telekom
D. Black
EMC Corporation
November 12, 2014

DiffServ interconnection classes and practice
draft-geib-tsvwg-diffserv-intercon-08

Abstract

This document proposes a limited and well defined set of DiffServ PHBs and codepoints to be applied at (inter)connections of two separately administered and operated networks. Many network providers operate MPLS using Treatment Aggregates for traffic marked with different DiffServ PHBs, and use MPLS for interconnection with other networks. This document offers a simple interconnection approach that may simplify operation of DiffServ for network interconnection among providers.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on May 16, 2015.

Copyright Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect

to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
1.1. Related work	4
2. MPLS and the Short Pipe tunnel model	5
3. An Interconnection class and codepoint scheme	6
3.1. End-to-end QoS: PHB and DS CodePoint Transparency	11
3.2. Treatment of Network Control traffic at carrier interconnection interfaces	12
4. Acknowledgements	13
5. IANA Considerations	13
6. Security Considerations	13
7. References	13
7.1. Normative References	13
7.2. Informative References	14
Appendix A. Annex A Carrier interconnection related DiffServ aspects	15
Appendix B. Annex 2 The MPLS Short Pipe Model and IP traffic	17
Appendix C. Change log	21
Authors' Addresses	21

1. Introduction

DiffServ has been deployed in many networks. As described by section 2.3.4.2 of RFC 2475, remarking of packets at domain boundaries is a DiffServ feature [RFC2475]. This draft proposes a set of standard QoS classes and code points at interconnection points to which and from which locally used classes and code points should be mapped.

RFC2474 specifies the DiffServ Codepoint Field [RFC2474]. Differentiated treatment is based on the specific DSCP. Once set, it may change. If traffic marked with unknown or unexpected DSCPs is received, RFC2474 recommends forwarding that traffic with default (best effort) treatment without changing the DSCP markings. Many networks do not follow this recommendation, and instead remark unknown or unexpected DSCPs to the zero DSCP for consistency with default (best effort) forwarding.

Many providers operate MPLS-based backbones that employ backbone traffic engineering to ensure that if a major link, switch, or router fails, the result will be a routed network that continues to meet its Service Level Agreements (SLAs). Based on that foundation, foundation, [RFC5127] introduces the concept of DiffServ Treatment

Aggregates, which enable traffic marked with multiple DSCPs to be forwarded in a single MPLS Traffic Class (TC). Like RFC 5127, this document assumes robust provider backbone traffic engineering.

RFC5127 recommends transmission of DSCPs as they are received. This is not possible, if the receiving and the transmitting domains at a network interconnection use different DSCPs for the PHBs involved.

This document is motivated by requirements for IP network interconnection with DiffServ support among providers that operate MPLS in their backbones, but is applicable to other technologies. The operational simplifications and methods in this document help align IP DiffServ functionality with MPLS limitations, particularly when MPLS penultimate hop popping is used. That is an important reason why this document specifies 4 interconnection Treatment Aggregates. Limiting DiffServ to a small number Treatment Aggregates can help ensure that network traffic leaves a network with the same DSCPs that it was received with. The approach proposed here may be extended by operators or future specifications.

In isolation, use of standard interconnection PHBs and DSCPs may appear to be additional effort for a network operator. The primary offsetting benefit is that the mapping from or to the interconnection PHBs and DSCPs is specified once for all of the interconnections to other networks that can use this approach. Otherwise, the PHBs and DSCPs have to be negotiated and configured independently for each network interconnection, which has poor scaling properties. Further, end-to-end QoS treatment is more likely to result when an interconnection code point scheme is used because traffic is remarked to the same PHBs at all network interconnections. This document supports one-to-one DSCP remarking at network interconnections (not n DSCP to one DSCP remarking).

The example given in RFC 5127 on aggregation of DiffServ service classes uses 4 Treatment Aggregates, and this document does likewise because:

- o The available coding space for carrying QoS information (e.g., DiffServ PHB) in MPLS and Ethernet is only 3 bits in size, and is intended for more than just QoS purposes (see e.g. [RFC5129]).
- o There should be unused codes for interconnection purposes. This leaves space for future standards, for private bilateral agreements and for local use PHBs and DSCPs.
- o Migrations from one code point scheme to another may require spare QoS code points.

RFC5127 provides recommendations on aggregation of DSCP-marked traffic into MPLS Treatment Aggregates and offers a deployment example [RFC5127] that does not work for the MPLS Short Pipe model when that model is used for ordinary network traffic. This document supports the MPLS Short Pipe model for ordinary network traffic and hence differs from the RFC5127 approach as follows:

- o remarking of received DSCPs to domain internal DSCPs is to be expected for ordinary IP traffic at provider edges (and for outer headers of tunneled IP traffic).
- o document follows RFC4594 in the proposed marking of provider Network Control traffic and expands RFC4594 on treatment of CS6 marked traffic at interconnection points (see section 3.2).

This document is organized as follows: section 2 reviews the MPLS Short Pipe tunnel model for DiffServ Tunnels [RFC3270]; effective support for that model is a crucial goal of this document. Section 3 introduces DiffServ interconnection Treatment Aggregates, plus the PHBs and DSCPs that are mapped to these Treatment Aggregates. Further, section 3 discusses treatment of non-tunneled and tunneled IP traffic and MPLS VPN QoS aspects. Finally Network Management PHB treatment is described. Annex A discusses how domain internal IP layer QoS schemes impact interconnection. Annex B describes the impact of the MPLS Short Pipe model (pen ultimate hop popping) on QoS related IP interconnections.

1.1. Related work

In addition to the activities that triggered this work, there are additional RFCs and Internet-drafts that may benefit from an interconnection PHB and DSCP scheme. RFC 5160 suggests Meta-QoS-Classes to enable deployment of standardized end to end QoS classes [RFC5160]. In private discussion, the authors of that RFC agree that the proposed interconnection class- and codepoint scheme and its enablement of standardised end to end classes would complement their own work.

Work on signaling Class of Service at interconnection interfaces by BGP [I-D.knoll-idr-cos-interconnect], [ID.idr-sla] is beyond the scope of this draft. When the basic DiffServ elements for network interconnection are used as described in this document, signaled access to QoS classes may be of interest. These two BGP documents focus on exchanging SLA and traffic conditioning parameters and assume that common PHBs identified by the signaled DSCPs have been established prior to BGP signaling of QoS.

2. MPLS and the Short Pipe tunnel model

The Pipe and Uniform models for Differentiated Services and Tunnels are defined in [RFC2983]. RFC3270 adds the MPLS Short Pipe model in order to support penultimate hop popping (PHP) of MPLS Labels, primarily for IP tunnels and VPNs. The Short Pipe model and PHP have become popular with many network providers that operate MPLS networks and are now widely used for ordinary network traffic, not just traffic encapsulated in IP tunnels and VPNs. This has important implications for DiffServ functionality in MPLS networks.

RFC 2474's recommendation to forward traffic with unrecognized DSCPs with Default (best effort) service without rewriting the DSCP has proven to be a poor operational practice. Network operation and management are simplified when there is a 1-1 match between the DSCP marked on the packet and the forwarding treatment (PHB) applied by network nodes. When this is done, CS0 (the all-zero DSCP) is the only DSCP used for Default forwarding of best effort traffic, so a common practice is to use CS0 to remark traffic received with unrecognized or unsupported DSCPs at network edges.

MPLS networks are more subtle in this regard, as it is possible to encode the provider's DSCP in the MPLS TC field and allow that to differ from the PHB indicated by the DSCP in the MPLS-encapsulated IP packet. That would allow an unrecognized DSCP to be carried edge-to-edge over an MPLS network, because the effective DSCP used by the MPLS network would be encoded in the MPLS label TC field (and also carried edge-to-edge); this approach assumes that a provider MPLS label with the provider's TC field being present at all hops within the provider's network.

The Short Pipe tunnel model and PHP violate that assumption because PHP pops and discards the MPLS provider label carrying the provider's TC field. That discard occurs one hop upstream of the MPLS tunnel endpoint, resulting in no provider TC info being available at tunnel egress. Therefore the DSCP field in the MPLS-encapsulated IP header has to contain a DSCP that is valid for the provider's network; propagating another DSCP edge-to-edge requires an IP tunnel of some form. In the absence of IP tunneling (a common case for MPLS networks), it is not possible to pass all 64 possible DSCP values edge-to-edge across an MPLS network. See Annex B for a more detailed discussion.

If transport of a large number (much greater than 4) DSCPs is required across a network that supports this DiffServ interconnection scheme, a tunnel or VPN can be provisioned for this purpose, so that the inner IP header carries the DSCP that is to be preserved not to be changed. From a network operations perspective, the customer

equipment (CE) is the preferred location for tunnel termination, although a receiving domains Provider Edge router is another viable option.

3. An Interconnection class and codepoint scheme

At an interconnection, the networks involved need to agree on the PHBs used for interconnection and the specific DSCP for each PHB. This may involve remarking for the interconnection; such remarking is part of the DiffServ Architecture [RFC2475], at least for the network edge nodes involved in interconnection. See Annex A for a more detailed discussion. This draft proposes a standard interconnection set of 4 Treatment Aggregates with well-defined DSCPs to be aggregated by them. A sending party remarks DSCPs from internal schemes to the interconnection code points. The receiving party remarks DSCPs to her internal scheme. The set of DSCPs and PHBs supported across the two interconnected domains and the treatment of PHBs and DSCPs not recognized by the receiving domain should be part of the interconnect SLA.

RFC 5127's four treatment aggregates include a Network Control aggregate for routing protocols and OAM traffic that is essential for network operation administration, control and management. Using this aggregate as one of the four in RFC 5127 implicitly assumes that network control traffic is forwarded in potential competition with all other network traffic, and hence DiffServ must favor such traffic (e.g., via use of the CS6 codepoint) for network stability. That is a reasonable assumption for IP-based networks where routing and OAM protocols are mixed with all other types of network traffic; corporate networks are an example.

In contrast, mixing of all traffic is not a reasonable assumption for MPLS-based provider or carrier networks, where customer traffic is usually segregated from network control (routing and OAM) traffic via other means, e.g., network control traffic use of separate LSPs that can be prioritized over customer LSPs (e.g., for VPN service) via other means. This sort of network control traffic from customer traffic is also used for MPLS-based network interconnections. In addition, many customers of a network provider do not exchange Network Control traffic (e.g., routing) with the network provider. For these reasons, a separate Network Control traffic aggregate is not important for MPLS-based carrier or provider networks; when such traffic is not segregated from other traffic, it may reasonably share the Assured Elastic treatment aggregate (as RFC 5127 suggests for a situation in which only three treatment aggregates are supported).

In contrast, VoIP is emerging as a valuable and important class of network traffic for which network-provided QoS is crucial, as even

minor glitches are immediately apparent to the humans involved in the conversation.

For these reasons, the Diffserv Interconnect scheme in this document departs from the approach in RFC 5127 by not providing a Network Control traffic aggregate, and instead dedicating the fourth traffic aggregate for VoIP traffic. Network Control traffic may still be exchanged across network interconnections, see Section 3.2 for further discussion.

Similar approaches to use of a small number of traffic aggregates (including recognition of the importance of VoIP traffic) have been taken in related standards and recommendations from outside the IETF, e.g., Y.1566 [Y.1566], GSM IR.34 [IR.34] and MEF23.1 [MEF23.1].

The list of the four Diffserv Interconnect traffic aggregates follows, highlighting differences from RFC 5127 and the specific traffic classes from RFC 4594 that each class aggregates.

Telephony Service Treatment Aggregate: PHB EF, DSCP 101 110 and VOICE-ADMIT, DSCP 101100, see [RFC3246] , [RFC4594][RFC5865]. This Treatment Aggregate corresponds to RFC 5127's real time Treatment Aggregate definition regarding the queuing, but it is restricted to transport Telephony Service Class traffic in the sense of RFC 4594.

Bulk Real-Time Treatment Aggregate: This Treatment Aggregate is designed to transport PHB AF41, DSCP 100 010 (the other AF4 PHB group PHBs and DSCPs may be used for future extension of the set of DSCPs carried by this Treatment Aggregate). This Treatment Aggregate is designed to transport the portions of RFC 5127's Real Time Treatment Aggregate, which consume large amounts of bandwidth, namely Broadcast Video, Real-Time Interactive and Multimedia Conferencing. The treatment aggregate should be configured with a rate queue (which is in line with RFC 4594 for the mentioned traffic classes). As compared to RFC 5127, the number of DSCPs has been reduced to one (initially) and the proposed queuing mechanism. The latter is however in line with RFC4594.

Assured Elastic Treatment Aggregate This Treatment Aggregate consists of the entire AF3 PHB group AF3, i.e., DSCPs 011 010, 011 100 and 011 110. As compared to RFC5127, just the number of DSCPs, which has been reduced. This document suggests to transport signaling marked by AF31. RFC5127 suggests to map Network Management traffic into this Treatment Aggregate, if no separate Network Control Treatment Aggregate is supported (for a more detailed discussion of

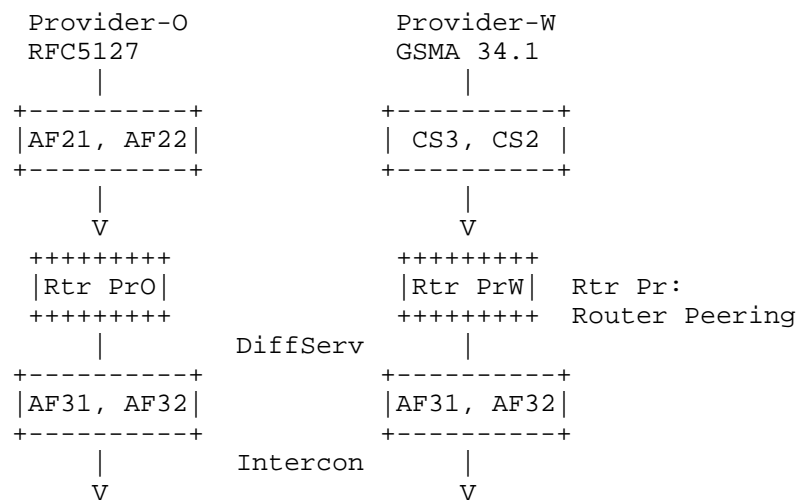
Network Control PHB treatment see section 3.2). GSMA IR.34 proposes to transport signaling traffic by AF31 too.

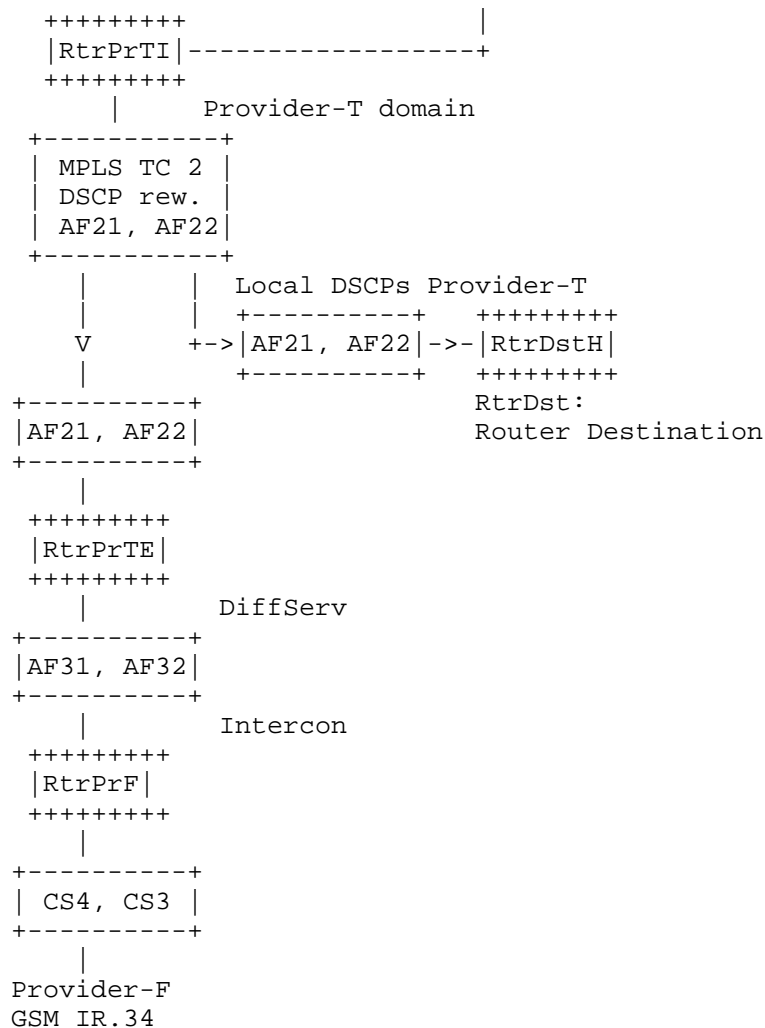
Default / Elastic Treatment Aggregate: transports the default PHB, CS0 with DSCP 000 000. RFC 5127 example refers to this Treatment Aggregate as Aggregate Elastic. An important difference as compared to RFC5127 is that any traffic with unrecognized or unsupported DSCPs may be remarked to this DSCP.

RFC 4594's Multimedia Streaming class has not been mapped to the above scheme. By the time of writing, the most popular streaming applications use TCP transport and adapt picture quality in the case of congestion. These applications are proprietary and still change behaviour frequently. At this state, the Bulk Real-Time Treatment Aggregate or the Bulk Real-Time Treatment Aggregate may be a reasonable match.

The overall approach to DSCP marking at network interconnections is illustrated by the following example. Provider O and provider W are peered with provider T. They have agreed upon a QoS interconnection SLA.

Traffic of provider O terminates within provider Ts network, while provider W's traffic transits through the network of provider T to provider F. Assume all providers to run their own internal codepoint schemes for a PHB group with properties of the DiffServ Intercon Assured Treatment Aggregate.





DiffServ Intercon example

Figure 1

It is easily visible that all providers only need to deploy internal DSCP to DiffServ Intercon DSCP mappings to exchange traffic in the desired classes. Provider W has decided that the properties of his internal classes CS3 and CS2 are best met by the Diffserv Intercon Assured Elastic Treatment Aggregate, PHBs AF31 and AF32 respectively. At the outgoing peering interface connecting provider W with provider

T remarks CS3 traffic to AF31 and CS 2 traffic to CS32. The domain internal PHBs of provider T meeting the Diffserv Intercon Assured Elastic Treatment Aggregate requirements is AF2. Hence AF31 traffic received at the interconnection with provider T is remarked to AF21 by the peering router of domain T. As domain T deploys MPLS, further the MPLS TC is set to 2. Traffic received with AF32 is remarked to AF22. The MPLS TC of the Treatment Aggregate is the same, TC 2. At the pen-ultimate MPLS node, the top MPLS label is removed. The packet should be forwarded as determined by the incoming MPLS TC. The peering router connecting domain T with domain F classifies the packet by its domain T internal DSCP AF21 for the Diffserv Intercon Assured Elastic Treatment Aggregate. As it leaves domain T on the interface to domain F, it is remarked to AF31. The peering router of domain F classifies the packet for domain F internal PHB CS4, as this is the PHB with properties matching Diffserv Intercon's Assured Elastic Treatment Aggregate. Likewise, AF21 traffic is remarked to AF32 by the peering router of domain T when leaving it and from AF32 to CS3 by domain F's peering router when receiving it.

This example can be extended. Suppose Provider-O also supports a PHB marked by CS2 and this PHB is supposed to be transported by QoS within Provider-T domain. Then Provider-O will remark it with a DSCP other than AF31 DSCP in order to preserve the differentiation from CS2; AF11 is one possibility that might be private to the interconnection between Provider-O and Provider-T; there's no assumption that Provider-W can also use AF11, as it may not be in the SLA with Provider-W.

Now suppose Provider-W supports CS2 for internal use only. Then no Diffserv intercon DSCP mapping may be configured at the peering router. Traffic, sent by Provider-W to Provider-T marked by CS2 due to a misconfiguration may be remarked to CS0 by Provider-T.

See section 3.1 for further discussion of this and DSCP transparency in general.

RFC5127 specifies a separate Treatment Aggregate for network control traffic. It may be present at interconnection interfaces too, but depending on the agreement between providers, Network Control traffic may also be classified into a different interconnection class. See section 3.2 for a detailed discussion on the treatment of Network Control traffic.

RFC2575 states that Ingress nodes must condition all other inbound traffic to ensure that the DS codepoints are acceptable; packets found to have unacceptable codepoints must either be discarded or must have their DS codepoints modified to acceptable values before being forwarded. For example, an ingress node receiving traffic from

a domain with which no enhanced service agreement exists may reset the DS codepoint to the Default PHB codepoint. As a consequence, an interconnect SLA needs to specify not only the treatment of traffic that arrives with a supported interconnect DSCP, but also the treatment of traffic that arrives with unsupported or unexpected DSCPs.

The proposed interconnect class and code point scheme is designed for point to point IP layer interconnections among MPLS networks. Other types of interconnections are out of scope of this document. The basic class and code point scheme is applicable on Ethernet layer too, if a provider e.g. supports Ethernet priorities like specified by IEEE 802.1p.

3.1. End-to-end QoS: PHB and DS CodePoint Transparency

This section describes how the use of a common PHB and DSCP scheme for interconnection can lead to end-to-end DiffServ-based QoS across networks that do not have common policies or practices for PHB and DSCP usage. This will initially be possible for PHBs and DSCPs corresponding to at most 3 or 4 Treatment Aggregates due to the MPLS considerations discussed previously.

Networks can be expected to differ in the number of PHBs available at interconnections (for terminating or transit service) and the DSCP values used within their domain. At an interconnection, Treatment Aggregate and PHB properties are best described by SLAs and related explanatory material. See annex A for a more detailed discussion about why PHB and DSCP usage is likely to differ among networks. For the above reasons and the desire to support interconnection among networks with different DiffServ schemes, the DiffServ interconnection scheme supports a small number of PHBs and DSCPs; this scheme is expandable.

The basic idea is that traffic sent with a DiffServ interconnect PHB and DSCP is restored to that PHB and DSCP (or a PHB and DSCP within the AF3 PHB group for the Assured Treatment Aggregate) at each network interconnection, even though a different PHB and DSCP may be used by each network involved. So, Bulk Inelastic traffic could be sent with AF41, remarked to CS3 by the first network and back to AF41 at the interconnection with the second network, which could mark it to CS5 and back to AF41 at the next interconnection, etc. The result is end-to-end QoS treatment consistent with the Bulk Inelastic Traffic Aggregate, and that is signaled or requested by the AF41 DSCP at each network interconnection in a fashion that allows each network operator to use their own internal PHB and DSCP scheme.

The key requirement is that the network ingress interconnect DSCP be restored at network egress, and a key observation is that this is only feasible in general for a small number of DSCPs.

3.2. Treatment of Network Control traffic at carrier interconnection interfaces

As specified by RFC4594, section 3.2, Network Control (NC) traffic marked by CS6 is to be expected at interconnection interfaces. This document does not change NC specifications of RFC4594, but observes that network control traffic received at network ingress is generally different from network control traffic within a network that is the primary use of CS6 envisioned by RFC 4594. A specific example is that some CS6 traffic exchanged across carrier interconnections is terminated at the network ingress node (e.g., if BGP is running between two routers on opposite ends of an interconnection link), which is consistent with RFC 4594's recommendation to not use CS6 when forwarding CS6-marked traffic originating from user-controlled end points.

The end-to-end QoS discussion in the previous section (3.1) is generally inapplicable to network control traffic - network control traffic is generally intended to control a network, not be transported across it. One exception is that network control traffic makes sense for a purchased transit agreement, and preservation of CS6 for network control traffic that is transited is reasonable in some cases. Use of an IP tunnel is suggested in order to reduce the risk of CS6 markings on transiting network control traffic being interpreted by the network providing the transit.

If the MPLS Short Pipe model is deployed for non tunneled IPv4 traffic, an IP network provider should limit access to the CS6 and CS7 DSCPs so that they are only used for network control traffic for the provider's own network.

Interconnecting carriers should specify treatment of CS6 marked traffic received at a carrier interconnection which is to be forwarded beyond the ingress node. An SLA covering the following cases is recommended when a provider wishes to send CS6 marked traffic across an interconnection link which isn't terminating at the interconnected ingress node:

- o classification of traffic which is network control traffic for both domains. This traffic should be classified and marked for the NC PHB.
- o classification of traffic which is network control traffic for the sending domain only. This traffic should be classified for a PHB

offering similar properties as the NC class (e.g. AF31 as specified by this document). As an example GSMA IR.34 proposes an Interactive class / AF31 to carry SIP and DIAMETER traffic. While this is service control traffic of high importance to the interconnected Mobile Network Operators, it is certainly no Network Control traffic for a fixed network providing transit. The example may not be perfect. It was picked nevertheless because it refers to an existing standard.

- o any other CS6 marked traffic should be remarked or dropped.

4. Acknowledgements

Al Morton and Sebastien Jobert provided feedback on many aspects during private discussions. Mohamed Boucadair and Thomas Knoll helped adding awareness of related work. Fred Baker and Brian Carpenter provided intensive feedback and discussion.

5. IANA Considerations

This memo includes no request to IANA.

6. Security Considerations

This document does not introduce new features, it describes how to use existing ones. The security section of RFC 2475 [RFC2475] and RFC 4594 [RFC4594] apply.

7. References

7.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, December 1998.
- [RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and W. Weiss, "An Architecture for Differentiated Services", RFC 2475, December 1998.
- [RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", RFC 2597, June 1999.

- [RFC3246] Davie, B., Charny, A., Bennet, J., Benson, K., Le Boudec, J., Courtney, W., Davari, S., Firoiu, V., and D. Stiliadis, "An Expedited Forwarding PHB (Per-Hop Behavior)", RFC 3246, March 2002.
- [RFC3260] Grossman, D., "New Terminology and Clarifications for Diffserv", RFC 3260, April 2002.
- [RFC3270] Le Faucheur, F., Wu, L., Davie, B., Davari, S., Vaananen, P., Krishnan, R., Cheval, P., and J. Heinanen, "Multi-Protocol Label Switching (MPLS) Support of Differentiated Services", RFC 3270, May 2002.
- [RFC5129] Davie, B., Briscoe, B., and J. Tay, "Explicit Congestion Marking in MPLS", RFC 5129, January 2008.
- [RFC5462] Andersson, L. and R. Asati, "Multiprotocol Label Switching (MPLS) Label Stack Entry: "EXP" Field Renamed to "Traffic Class" Field", RFC 5462, February 2009.
- [RFC5865] Baker, F., Polk, J., and M. Dolly, "A Differentiated Services Code Point (DSCP) for Capacity-Admitted Traffic", RFC 5865, May 2010.
- [min_ref] authSurName, authInitials., "Minimal Reference", 2006.

7.2. Informative References

- [I-D.knoll-idr-cos-interconnect]
Knoll, T., "BGP Class of Service Interconnection", draft-knoll-idr-cos-interconnect-13 (work in progress), November 2014.
- [ID.idr-sla]
IETF, "Inter-domain SLA Exchange", IETF,
<http://datatracker.ietf.org/doc/draft-ietf-idr-sla-exchange/>, 2013.
- [IEEE802.1Q]
IEEE, "IEEE Standard for Local and Metropolitan Area Networks - Virtual Bridged Local Area Networks", 2005.
- [IR.34]
GSMA Association, "IR.34 Inter-Service Provider IP Backbone Guidelines Version 7.0", GSMA, GSMA IR.34
<http://www.gsma.com/newsroom/wp-content/uploads/2012/03/ir.34.pdf>, 2012.

- [MEF23.1] MEF, "Implementation Agreement MEF 23.1 Carrier Ethernet Class of Service Phase 2", MEF, MEF23.1
http://metroethernetforum.org/PDF_Documents/technical-specifications/MEF_23.1.pdf, 2012.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", RFC 2983, October 2000.
- [RFC4594] Babiarz, J., Chan, K., and F. Baker, "Configuration Guidelines for DiffServ Service Classes", RFC 4594, August 2006.
- [RFC5127] Chan, K., Babiarz, J., and F. Baker, "Aggregation of Diffserv Service Classes", RFC 5127, February 2008.
- [RFC5160] Levis, P. and M. Boucadair, "Considerations of Provider-to-Provider Agreements for Internet-Scale Quality of Service (QoS)", RFC 5160, March 2008.
- [Y.1566] ITU-T, "Quality of service mapping and interconnection between Ethernet, IP and multiprotocol label switching networks", ITU,
<http://www.itu.int/rec/T-REC-Y.1566-201207-I/en>, 2012.

Appendix A. Annex A Carrier interconnection related DiffServ aspects

This annex provides a general discussion of PHB and DSCP mapping at IP interconnection interfaces. It also informs about limitations and likely DSCP changes.

The following scenarios start from a domain sending non-tunneled IP traffic using a PHB and a corresponding DSCP to an interconnected domain. The receiving domain may

- o Support the PHB and offer the same corresponding DSCP.
- o Not support the PHB and use the DSCP for a different PHB.
- o Not support the PHB and not use the DSCP.
- o Support the PHB with a differing DSCP, and the DSCP of the sending domain is not used for another PHB
- o Support the PHB with a differing DSCP, and the DSCP of the sending domain is used for another PHB.

RFC2475 allows for local use PHB groups which are only available within a domain. If such a local use PHB is present, non-tunneled IP traffic possibly cannot utilize 64 DSCPs end-to-end.

If a domain receives traffic for a PHB, which it does not support, there are two general scenarios:

- o The received DSCP is not available for usage within the domain.
- o The received DSCP is available for usage within the domain.

RFC2474 suggests to transport packets received with unrecognized DSCPs by the default PHB and leave the DSCP as received. Also if a particular DSCP is spare within a domain, it may later change its QoS design and assign a PHB to a formerly unused DSCP (which a customer used to transit through this unrecognized DSCP will note, as his DSCP will be remarked). A transparent transport of the same DSCP as unknown with the default PHB may no longer be possible. Remarking to another DSCP apart from the Default PHBs DSCP does not seem to be a good option in the latter case. Which other DSCP is making sense? If a domain interconnects with many other domains, the questions asked here may have to be answered multiple times.

The scenarios above indicate, that reliably delivering a non-tunneled IP packet by the same PHB and DSCP unchanged end-to-end is only likely, if both domains support this DSCP and use the same corresponding DSCP.

Limitations in the number of supported PHBs are to be expected if DiffServ is applied across different domains. Unchanged end-to-end DSCPs should only be expected for non-tunneled IP traffic, if the PHB and DSCP are well specified and generally deployed. This is true for Default Forwarding. EF PHB is a candidate. The Network Control PHB is a local use only example, hence end-to-end support of CS6 for non-tunneled IP traffic at interconnection points should only be expected, if the receiving domain regards this traffic as Network Control traffic relevant for the own domain too.

DiffServ Intercon proposes a well defined set of PHBs and corresponding DSCPs at interconnection points. A PHB to DSCPs correspondence is specified at least for interconnection interfaces. Supported PHBs should be available end-to-end, but domain internal DSCPs may change end-to-end.

Appendix B. Annex 2 The MPLS Short Pipe Model and IP traffic

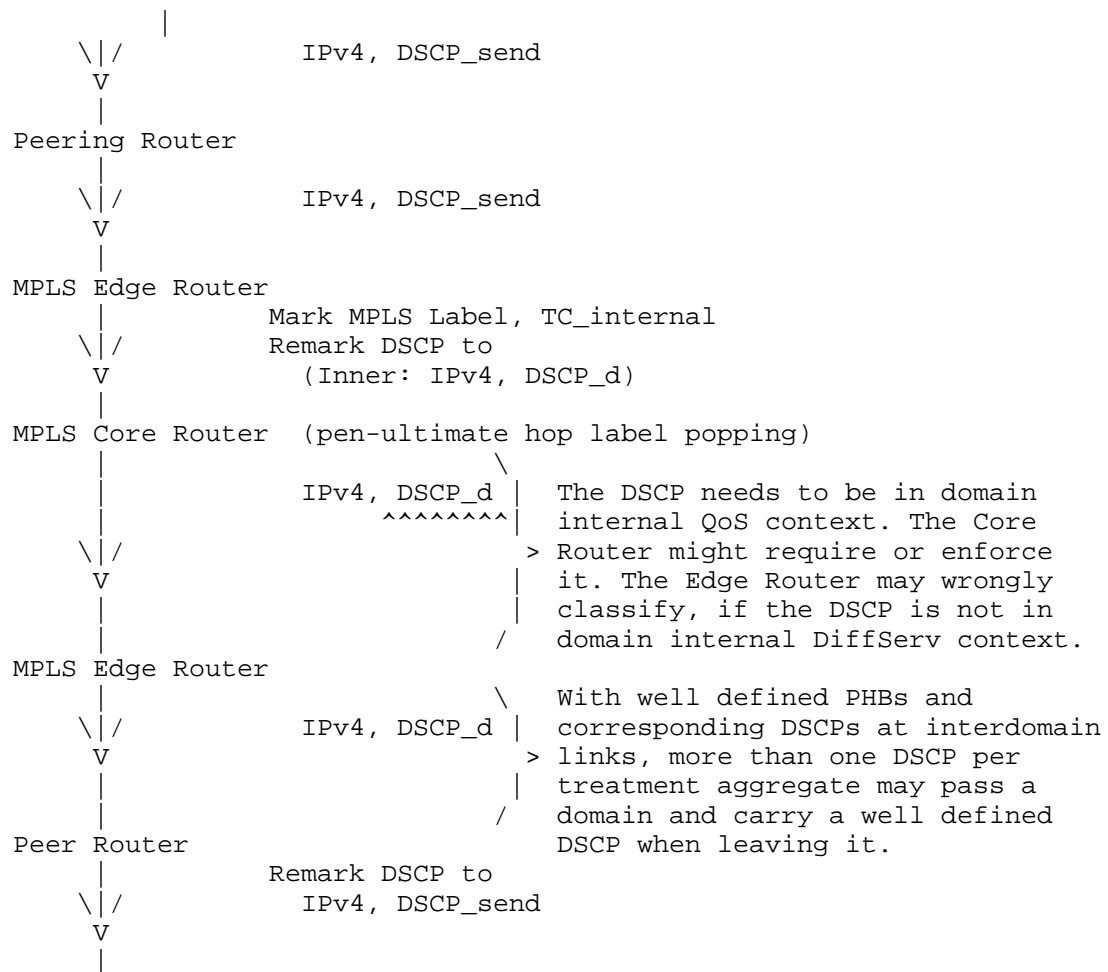
The MPLS Short Pipe Model (or Pen-ultimate Hop Label Popping) is widely deployed by IP carriers. If non-tunneled IPv4 traffic is transported using MPLS Short Pipe, IP headers appear inside the last section of the MPLS domain. This likely impacts the number of PHBs and DSCPs a network provider supports for this kind of traffic. Figure 2 provides an example for the treatment of this kind of traffic.

In the case of tunneled IPv4 traffic, only the outer tunnel header is exposed. Assuming the tunnel not to terminate within the MPLS network section, only the outer tunnel DSCP is impacted.

Non-tunneled IPv6 traffic and Layer 2 and Layer 3 VPN traffic all use an additional label. Hence no IP header is exposed within an MPLS domain.

Carriers may first design their own QoS PHB and codepoint scheme before they worry about interconnection. PHB and corresponding codepoint schemes usually differ between different carriers. PHBs may be mapped. A DSCP rewrite should be expected at an interconnection interface at least for plain IP traffic.

RFC3270 suggests deployment of the Short Pipe Model only in the case of VPNs. State of the art deployments also support transport of non tunneled IPv4 traffic. This is shown in figure 2.



Short-Pipe / Pen-ultimate hop popping example

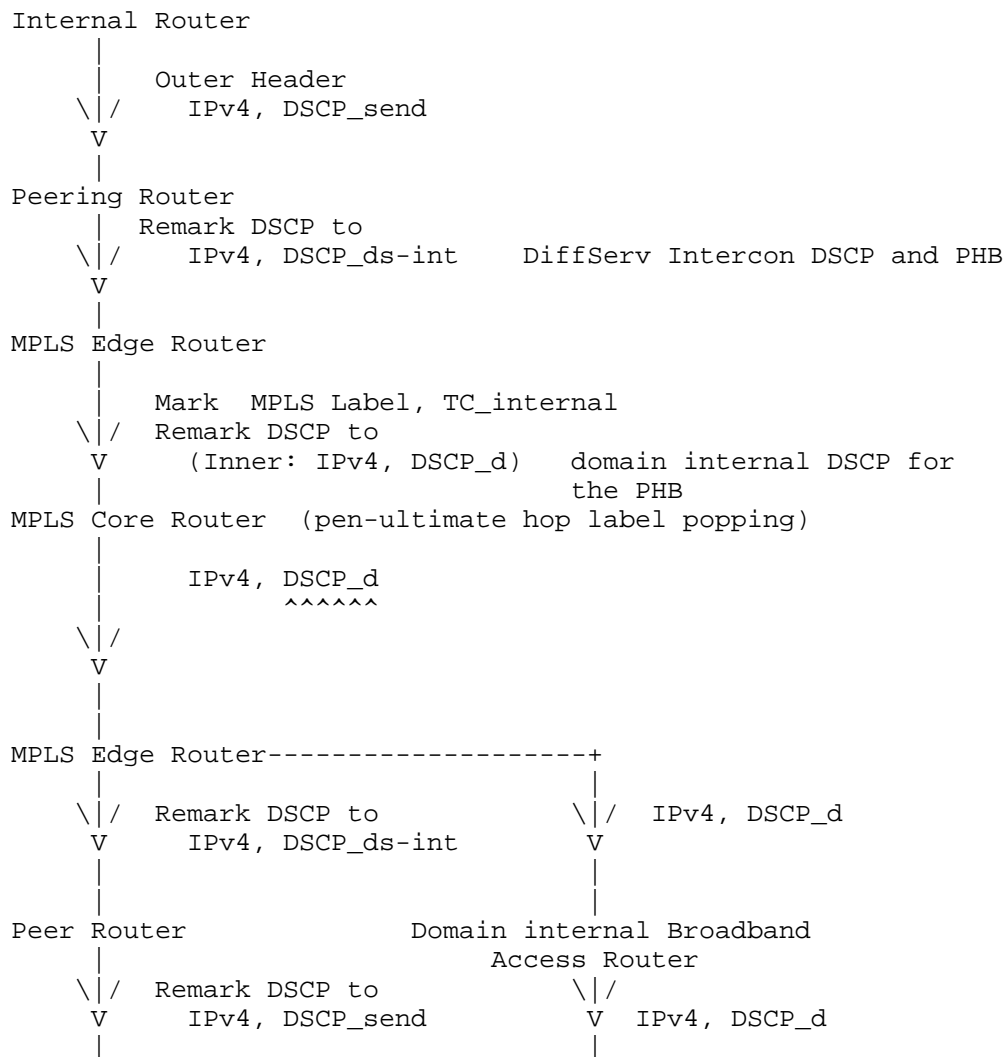
Figure 2

The packets IP DSCP must be in a well understood Diffserv context for schedulers and classifiers on the interfaces of the ultimate MPLS link. These are domain internal and a domain operating in this mode enforces DSCPs resulting in reliable domain internal QoS operation.

Without DiffServ-Intercon treatment, the traffic always leaves the domain having internal DS codepoints. DSCP_send of the figure above is remarked to the receiving domains DiffServ scheme. It leaves the

domain marked by the domains DSCP_d. Every carrier must deploy per peer PHB and DSCP mapping schemes.

If DiffServ-Intercon is applied, only traffic terminating within a domain must be aligned with the domain internal DiffServ Codepoint scheme. Traffic transiting through the domain can be easily mapped and remapped to an original DSCP. This is shown in figure 3. Of course the domain internal limitations caused by the Short Pipe model still apply.



Short-Pipe example with Diffserv-Intercon

Figure 3

Picking up terminology of RFC2983 and RFC3270, DiffServ intercon emulates the long pipe model for the PHBs it supports, if traffic is terminating in the receiving domain.

Looking at the peering interfaces only, for transiting QoS traffic DiffServ-Intercon emulates the uniform model for the PHBs and DSCPs

supported. Packets are expected to leave a domain with the DSCP/PHB as received (and per flow within each PHB in the same order as received). MPLS Treatment Aggregates should not experience congestion under standard operational conditions. The peering links need to be engineered to be congestion free too for QoS PHBs, if also the IP transit transport is to be congestion free.

Appendix C. Change log

- 00 to 01 Added terminology and references. Added details and information to interconnection class and codepoint scheme. Editorial changes.
- 01 to 02 Added some references regarding related work. Clarified class definitions. Further editorial improvements.
- 02 to 03 Consistent terminology. Discussion of Network Management PHB at interconnection interfaces. Editorial review.
- 03 to 04 Again improved terminology. Better wording of Network Control PHB at interconnection interfaces.
- 04 to 05 Large rewrite and re-ordering of contents.
- 05 to 06 Description of IP and MPLS related requirements and constraints on DSCP rewrites.
- 06 to 07 Largely rewrite, improved match and comparison with RFCs 4594 and 5127.
- 07 to 08 Added Annex A and B which were forgotten when putting together -07

Authors' Addresses

Ruediger Geib (editor)
Deutsche Telekom
Heinrich Hertz Str. 3-7
Darmstadt 64295
Germany

Phone: +49 6151 5812747
Email: Ruediger.Geib@telekom.de

David L. Black
EMC Corporation
176 South Street
Hopkinton, MA
USA

Phone: +1 (508) 293-7953
Email: david.black@emc.com

Network Working Group
Internet-Draft
Intended status: Best Current Practice
Expires: March 13, 2014

R. Stewart
Adara Networks
M. Tuexen
I. Ruengeler
Muenster Univ. of Appl. Sciences
September 09, 2013

Stream Control Transmission Protocol (SCTP) Network Address Translation
draft-ietf-behave-sctpnat-09.txt

Abstract

Stream Control Transmission Protocol [RFC4960] provides a reliable communications channel between two end-hosts in many ways similar to TCP [RFC0793]. With the widespread deployment of Network Address Translators (NAT), specialized code has been added to NAT for TCP that allows multiple hosts to reside behind a NAT and yet use only a single globally unique IPv4 address, even when two hosts (behind a NAT) choose the same port numbers for their connection. This additional code is sometimes classified as Network Address and Port Translation or NAPT. To date, specialized code for SCTP has NOT yet been added to most NATs so that only pure NAT is available. The end result of this is that only one SCTP capable host can be behind a NAT.

This document describes an SCTP specific variant of NAT which provides similar features of NAPT in the single point and multi-point traversal scenario.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on March 13, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	2
2. Conventions	3
3. Terminology	3
4. SCTP NAT Traversal Scenarios	4
4.1. Single Point Traversal	4
4.2. Multi Point Traversal	5
5. Limitations of Classical NAPT for SCTP	6
6. The SCTP Specific Variant of NAT	6
7. NAT to SCTP	10
8. Handling of Fragmented SCTP Packets	10
9. Various Examples of NAT Traversals	10
9.1. Single-homed Client to Single-homed Server	10
9.2. Single-homed Client to Multi-homed Server	12
9.3. Multihomed Client and Server	15
9.4. NAT Loses Its State	18
9.5. Peer-to-Peer Communication	20
10. IANA Considerations	24
11. Security Considerations	24
12. Acknowledgments	24
13. References	24
13.1. Normative References	24
13.2. Informative References	25
Authors' Addresses	25

1. Introduction

Stream Control Transmission Protocol [RFC4960] provides a reliable communications channel between two end-hosts in many ways similar to TCP [RFC0793]. With the widespread deployment of Network Address Translators (NAT), specialized code has been added to NAT for TCP that allows multiple hosts to reside behind a NAT and use private

addresses (see [RFC5735]) and yet use only a single globally unique IPv4 address, even when two hosts (behind a NAT) choose the same port numbers for their connection. This additional code is sometimes classified as Network Address and Port Translation or NAPT. To date, specialized code for SCTP has not yet been added to most NATs so that only true NAT is available. The end result of this is that only one SCTP capable host can be behind a NAT.

This document proposes an SCTP specific variant NAT that provides the NAPT functionality without changing SCTP port numbers. The authors feel it is possible and desirable to make these changes for a number of reasons.

- o It is desirable for SCTP internal end-hosts on multiple platforms to be able to share a NAT's public IP address, much as TCP does today.
- o If a NAT does not need to change any data within an SCTP packet it will reduce the processing burden of NAT'ing SCTP by NOT needing to execute the CRC32c checksum required by SCTP.
- o Not having to touch the IP payload makes the processing of ICMP messages in NATs easier.

2. 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. Terminology

For this discussion we will use several terms, which we will define and point out in Figure 1.

Private-Address (Priv-Addr): The private address that is known to the internal host.

Internal-Port (Int-Port): The port number that is in use by the host holding the Private-Address.

Internal-VTag (Int-VTag): The Verification Tag that the internal host has chosen for its communication. The VTag is a unique 32 bit tag that must accompany any incoming SCTP packet for this association to the Private-Address.

External-Address (Ext-Addr): The address that an internal host is attempting to contact.

External-Port (Ext-Port): The port number of the peer process at the External-Address.

External-VTag (Ext-VTag): The Verification Tag that the host holding the External-Address has chosen for its communication. The VTag is a unique 32 bit tag that must accompany any incoming SCTP packet for this association to the External-Address.

Public-Address (Pub-Addr): The public address assigned to the NAT box which it uses as a source address when sending packets towards the External-Address.

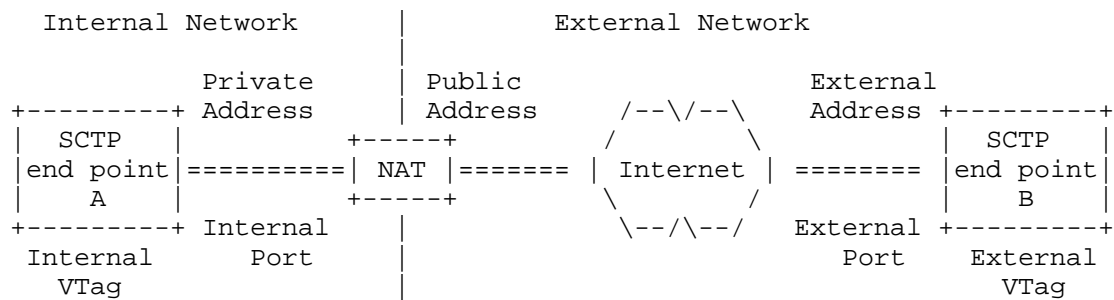


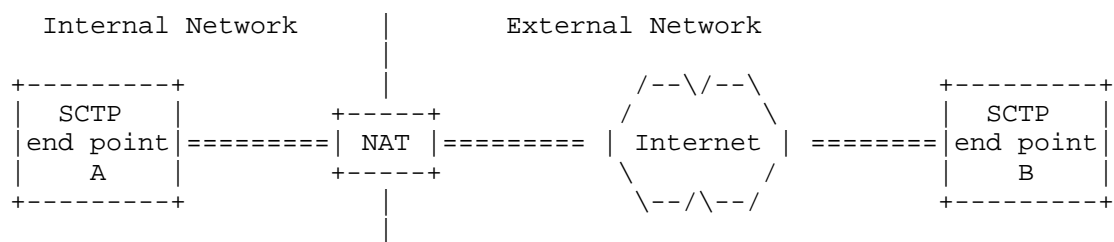
Figure 1: Architecture

4. SCTP NAT Traversal Scenarios

This section defines the notion of single and multi-point NAT traversal.

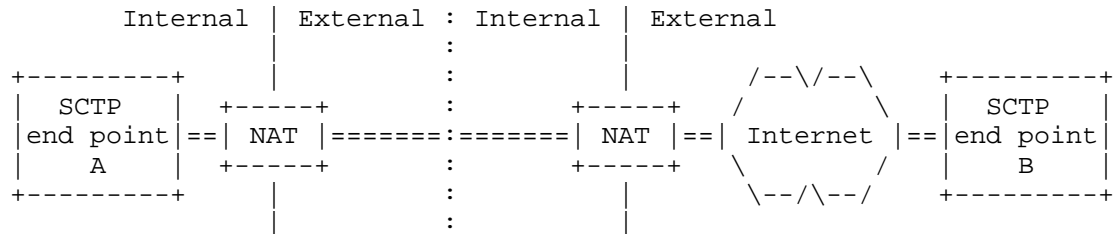
4.1. Single Point Traversal

In this case, all packets in the SCTP association go through a single NAT, as shown below:



Single NAT scenario

A variation of this case is shown below, i.e., multiple NATs in a single path:



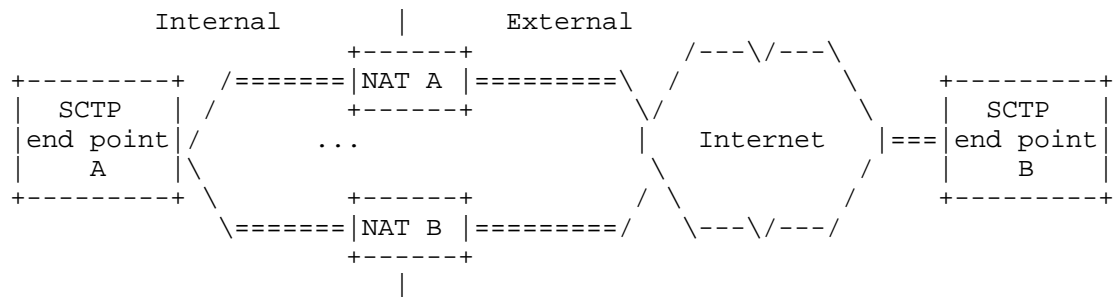
Serial NATs scenario

In this single point traversal scenario, we must acknowledge that while one of the main benefits of SCTP multi-homing is redundant paths, the NAT function represents a single point of failure in the path of the SCTP multi-home association. However, the rest of the path may still benefit from path diversity provided by SCTP multi-homing.

The two SCTP endpoints in this case can be either single-homed or multi-homed. However, the important thing is that the NAT (or NATs) in this case sees all the packets of the SCTP association.

4.2. Multi Point Traversal

This case involves multiple NATs and each NAT only sees some of the packets in the SCTP association. An example is shown below:



Parallel NATs scenario

This case does NOT apply to a single-homed SCTP association (i.e., BOTH endpoints in the association use only one IP address). The advantage here is that the existence of multiple NAT traversal points can preserve the path diversity of a multi-homed association for the

entire path. This in turn can improve the robustness of the communication.

5. Limitations of Classical NAT for SCTP

Using classical NAT may result in changing one of the SCTP port numbers during the processing which requires the recomputation of the transport layer checksum. Whereas for UDP and TCP this can be done very efficiently, for SCTP the checksum (CRC32c) over the entire packet needs to be recomputed. This would add considerable to the NAT computational burden, however hardware support may mitigate this in some implementations.

An SCTP endpoint may have multiple addresses but only has a single port number. To make multipoint traversal work, all the NATs involved must recognize the packets they see as belonging to the same SCTP association and perform port number translation in a consistent way. One possible way of doing this is to use pre-defined table of ports and addresses configured within each NAT. Other mechanisms could make use of NAT to NAT communication. Such mechanisms are considered by the authors not to be deployable on a wide scale base and thus not a recommended solution. Therefore the SCTP variant of NAT has been developed.

6. The SCTP Specific Variant of NAT

In this section we assume that we have multiple SCTP capable hosts behind a NAT which has one Public-Address. Furthermore we are focusing in this section on the single point traversal scenario.

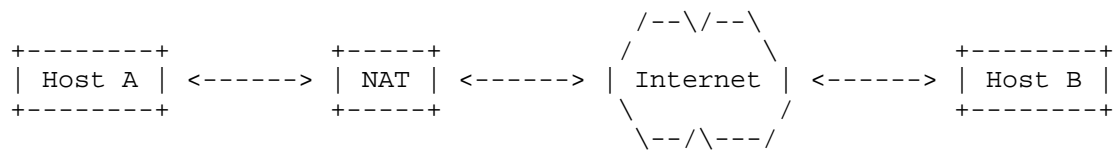
The modification of SCTP packets sent to the public Internet is easy. The source address of the packet has to be replaced with the Public-Address. It may also be necessary to establish some state in the NAT box to handle incoming packets, which is discussed later.

For SCTP packets coming from the public Internet the destination address of the packets has to be replaced with the Private-Address of the host the packet has to be delivered to. The lookup of the Private-Address is based on the External-VTag, External-Port, External-Address, Internal-VTag and the Internal-Port.

For the SCTP NAT processing the NAT box has to maintain a table of Internal-VTag, Internal-Port, Private-Address, External-VTag, External-Port and whether the restart procedure is disabled or not. An entry in that table is called a NAT state control block. The function Create() obtains the just mentioned parameters and returns a NAT-State control block.

The entries in this table fulfill some uniqueness conditions. There must not be more than one entry with the same pair of Internal-Port and External-Port. This rule can be relaxed, if all entries with the same Internal-Port and External-Port have the support for the restart procedure enabled. In this case there must be no more than one entry with the same Internal-Port, External-Port and Ext-VTag and no more than one entry with the same Internal-Port, External-Port and Int-VTag.

The processing of outgoing SCTP packets containing an INIT-chunk is described in the following figure. The scenario shown is valid for all message flows in this section.



```

                INIT[Initiate-Tag]
Priv-Addr:Int-Port -----> Ext-Addr:Ext-Port
                        Ext-VTag=0

                Create(Initiate-Tag, Int-Port, Priv-Addr, 0)
                Returns(NAT-State control block)

Translate To:

                INIT[Initiate-Tag]
Pub-Addr:Int-Port -----> Ext-Addr:Ext-Port
                        Ext-VTag=0

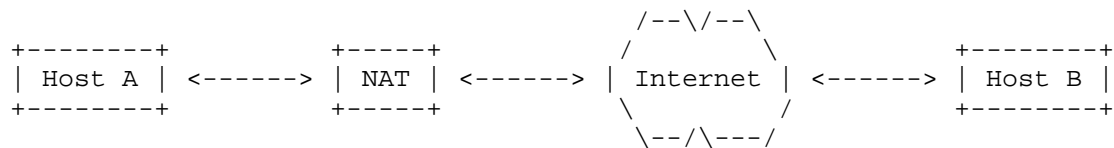
```

It should be noted that normally a NAT control block will be created. However, it is possible that there is already a NAT control block with the same External-Address, External-Port, Internal-Port, and Internal-VTag but different Private-Address. In this case the INIT SHOULD be dropped by the NAT and an ABORT SHOULD be sent back to the SCTP host with the M-Bit set and an appropriate error cause (see [I-D.ietf-tsvwg-natsupp] for the format). The source address of the packet containing the ABORT chunk MUST be the destination address of the packet containing the INIT chunk.

It is also possible that a connection to External-Address and External-Port exists without an Internal-VTag conflict but the

External-Address does not support the DISABLE_RESTART feature (noted in the NAT control block when the prior connection was established). In such a case the INIT SHOULD be dropped by the NAT and an ABORT SHOULD be sent back to the SCTP host with the M-Bit set and an appropriate error cause (see [I-D.ietf-tsvwg-natsupp] for the format).

The processing of outgoing SCTP packets containing no INIT-chunk is described in the following figure.

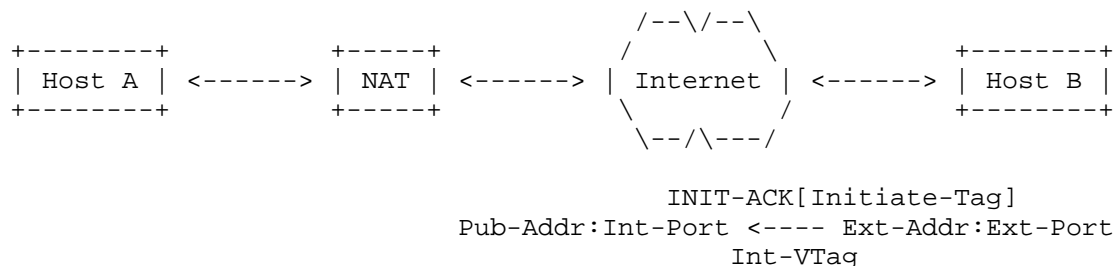


Priv-Addr:Int-Port -----> Ext-Addr:Ext-Port
 Ext-VTag

Translate To:

Pub-Addr:Int-Port -----> Ext-Addr:Ext-Port
 Ext-VTag

The processing of incoming SCTP packets containing INIT-ACK chunks is described in the following figure. The Lookup() function getting as input the Internal-VTag, Internal-Port, External-VTag (=0), External-Port, and External-Address, returns the corresponding entry of the NAT table and updates the External-VTag by substituting it with the value of the Initiate-Tag of the INIT-ACK chunk. The wildcard character signifies that the parameter's value is not considered in the Lookup() function or changed in the Update() function, respectively.



```

Lookup(Int-VTag, Int-Port, *, 0, Ext-Port)
Update(*, *, *, Initiate-Tag, *)

```

```

Returns(NAT-State control block containing Private-Address)

```

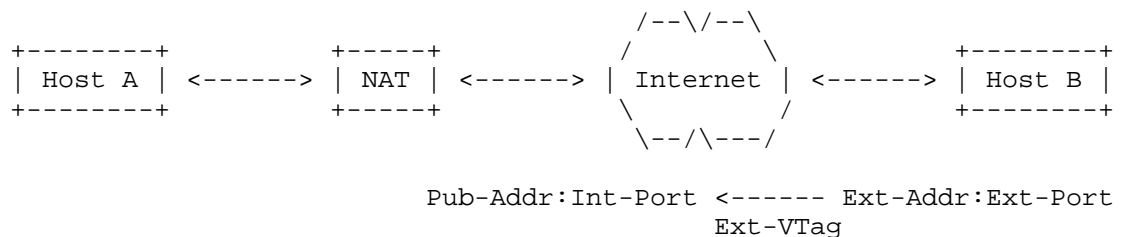
```

      INIT-ACK[Initiate-Tag]
Priv-Addr:Int-Port <----- Ext-Addr:Ext-Port
                    Int-VTag

```

In the case Lookup fails, the SCTP packet is dropped. The Update routine inserts the External-VTag (the Initiate-Tag of the INIT-ACK chunk) in the NAT state control block.

The processing of incoming SCTP packets containing an ABORT or SHUTDOWN-COMPLETE chunk with the T-Bit set is described in the following figure.



```

Lookup(0, Int-Port, *, Ext-VTag, Ext-Port)

```

```

Returns(NAT-State control block containing Private-Address)

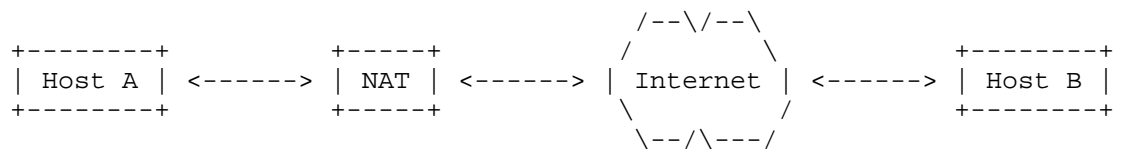
```

```

Priv-Addr:Int-Port <----- Ext-Addr:Ext-Port
                    Ext-VTag

```

The processing of other incoming SCTP packets is described in the following figure.



```

Pub-Addr: Int-Port <----- Ext-Addr: Ext-Port
                          Int-VTag

```

```

Lookup(Int-VTag, Int-Port, *, *, Ext-Port)

```

```

Returns(NAT-State control block containing Local-Address)

```

```

Priv-Addr: Int-Port <----- Ext-Addr: Ext-Port
                          Int-VTag

```

For an incoming packet containing an INIT-chunk a table lookup is made only based on the addresses and port numbers. If an entry with an External-VTag of zero is found, it is considered a match and the External-VTag is updated.

This allows the handling of INIT-collision through NAT.

7. NAT to SCTP

This document at various places discusses the sending of specialized SCTP chunks (e.g. an ABORT with M-Bit set). These chunks and procedures are not defined in this document, but instead are defined in [I-D.ietf-tsvwg-natsupp]. The NAT implementer should refer to [I-D.ietf-tsvwg-natsupp] for detailed descriptions of packet formats and procedures.

8. Handling of Fragmented SCTP Packets

A NAT box MUST support IP reassembly of received fragmented SCTP packets. The fragments may arrive in any order.

When an SCTP packet has to be fragmented by the NAT box and the IP header forbids fragmentation a corresponding ICMP packet SHOULD be sent.

9. Various Examples of NAT Traversals

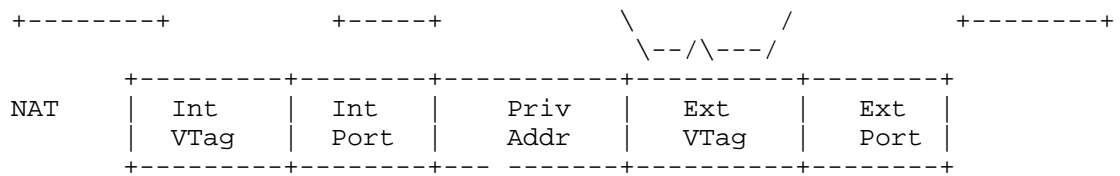
9.1. Single-homed Client to Single-homed Server

The internal client starts the association with the external server via a four-way-handshake. Host A starts by sending an INIT chunk.

```

+-----+          +-----+          /--\ /--\          +-----+
| Host A | <-----> | NAT   | <-----> | Internet | <-----> | Host B |

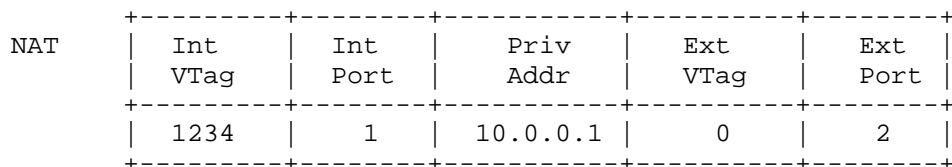
```



```
INIT[Initiate-Tag = 1234]
10.0.0.1:1 -----> 100.0.0.1:2
    Ext-VTtag = 0
```

A NAT entry is created, the source address is substituted and the packet is sent on:

NAT creates entry:

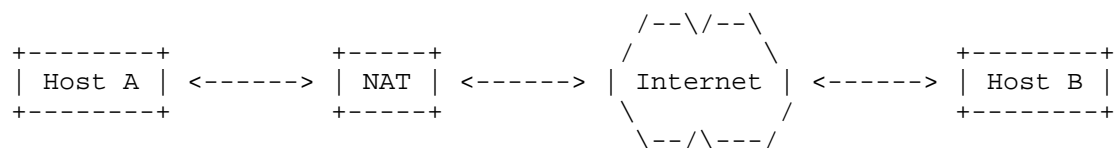


```

INIT[Initiate-Tag = 1234]
101.0.0.1:1 -----> 100.0.0.1:2
                     Ext-VTtag = 0

```

Host B receives the INIT and sends an INIT-ACK with the NAT's external address as destination address.

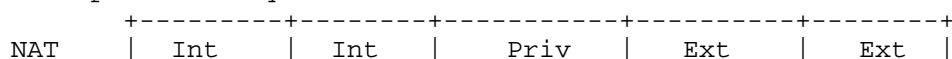


```

          INIT-ACK[Initiate-Tag = 5678]
101.0.0.1:1 <----- 100.0.0.1:2
                  Int-VTag = 1234

```

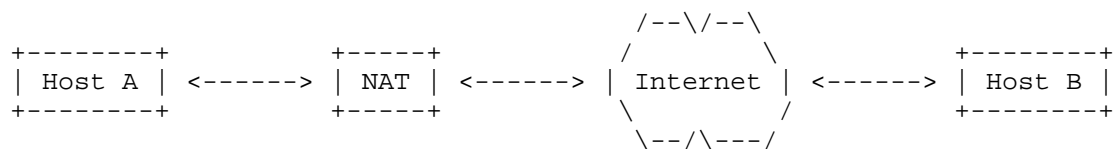
NAT updates entry:



VTag	Port	Addr	VTag	Port
1234	1	10.0.0.1	5678	2

```
INIT-ACK[Initiate-Tag = 5678]
10.0.0.1:1 <----- 100.0.0.1:2
      Int-VTag = 1234
```

The handshake finishes with a COOKIE-ECHO acknowledged by a COOKIE-ACK.



```
      COOKIE-ECHO
10.0.0.1:1 -----> 100.0.0.1:2
      Ext-VTag = 5678
```

```

                                COOKIE-ECHO
101.0.0.1:1 -----> 100.0.0.1:2
                                Ext-VTag = 5678
```

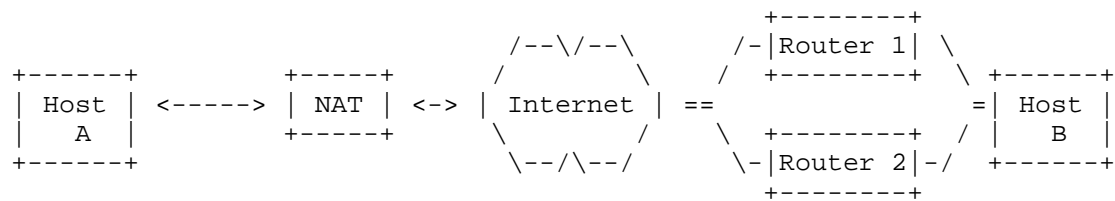
```

                                COOKIE-ACK
101.0.0.1:1 <----- 100.0.0.1:2
                                Int-VTag = 1234
```

```
      COOKIE-ACK
10.0.0.1:1 <----- 100.0.0.1:2
      Int-VTag = 1234
```

9.2. Single-homed Client to Multi-homed Server

The internal client is single-homed whereas the external server is multi-homed. The client (Host A) sends an INIT like in the single-homed case.



NAT	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port

```

INIT[Initiate-Tag = 1234]
10.0.0.1:1 ---> 100.0.0.1:2
      Ext-VTag = 0

```

NAT creates entry:

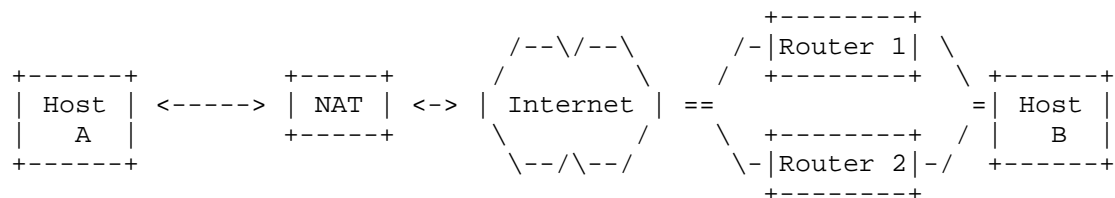
NAT	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
	1234	1	10.0.0.1	0	2

```

                                INIT[Initiate-Tag = 1234]
101.0.0.1:1 -----> 100.0.0.1:2
                        Ext-VTag = 0

```

The server (Host B) includes its two addresses in the INIT-ACK chunk, which results in two NAT entries.



```

      INIT-ACK[Initiate-tag = 5678, IP-Addr = 100.1.0.1]
101.0.0.1:1 <----- 100.0.0.1:2
                  Int-VTag = 1234

```

NAT does need to change the table for second address:

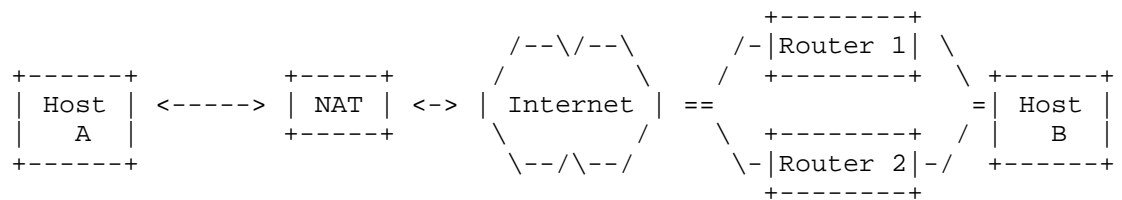
NAT						
	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port	
	1234	1	10.0.0.1	5678	2	

```

INIT-ACK[Initiate-Tag = 5678]
10.0.0.1:1 <--- 100.0.0.1:2
      Int-VTag = 1234

```

The handshake finishes with a COOKIE-ECHO acknowledged by a COOKIE-ACK.



```

      COOKIE-ECHO
10.0.0.1:1 ---> 100.0.0.1:2
      ExtVTag = 5678

```

```

                                     COOKIE-ECHO
101.0.0.1:1 -----> 100.0.0.1:2
                                     Ext-VTag = 5678

```

```

                                     COOKIE-ACK
101.0.0.1:1 <----- 100.0.0.1:2
                                     Int-VTag = 1234

```

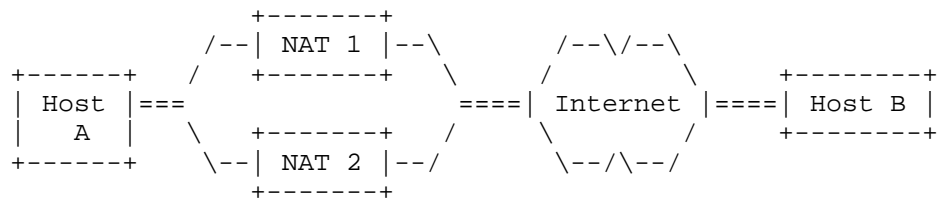
```

      COOKIE-ACK
10.0.0.1:1 <--- 100.0.0.1:2
      Int-VTag = 1234

```

9.3. Multihomed Client and Server

The client (Host A) sends an INIT to the server (Host B), but does not include the second address.



NAT 1	Int		Priv	Ext	
	VTag	Port		VTag	Port

```

      INIT[Initiate-Tag = 1234]
10.0.0.1:1 -----> 100.0.0.1:2
      Ext-VTag = 0

```

NAT 1 creates entry:

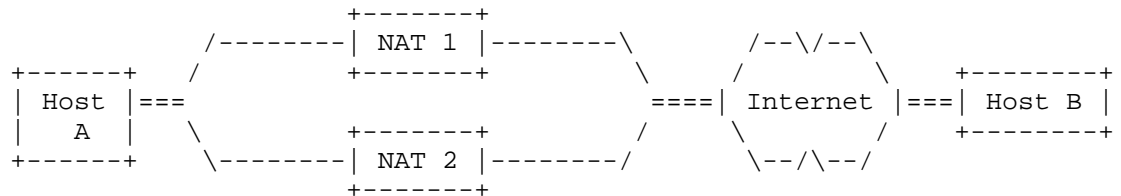
NAT 1	Int		Priv	Ext	
	VTag	Port		VTag	Port
	1234	1	10.0.0.1	0	2

```

      INIT[Initiate-Tag = 1234]
101.0.0.1:1 -----> 100.0.0.1:2
      ExtVTag = 0

```

Host B includes its second address in the INIT-ACK, which results in two NAT entries in NAT 1.



```

INIT-ACK[Initiate-Tag = 5678, IP-Addr = 100.1.0.1]
101.0.0.1:1 <----- 100.0.0.1:2
                    Int-VTag = 1234

```

NAT 1 does not need to update the table for second address:

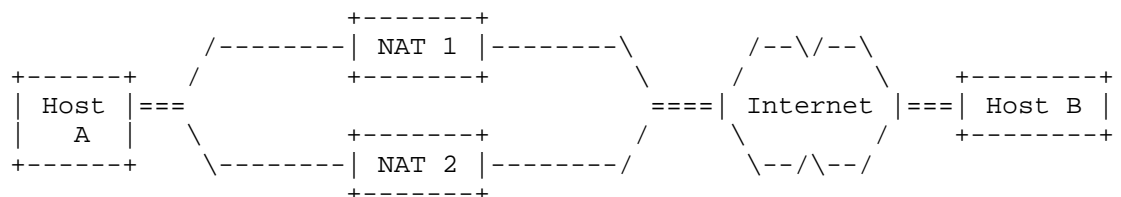
NAT 1					
	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
	1234	1	10.0.0.1	5678	2

```

INIT-ACK[Initiate-Tag = 5678]
10.0.0.1:1 <-----100.0.0.1:2
                    Int-VTag = 1234

```

The handshake finishes with a COOKIE-ECHO acknowledged by a COOKIE-ACK.



COOKIE-ECHO

```
10.0.0.1:1 -----> 100.0.0.1:2
      Ext-VTag = 5678
```

```

                                COOKIE-ECHO
101.0.0.1:1 -----> 100.0.0.1:2
                        Ext-VTag = 5678
```

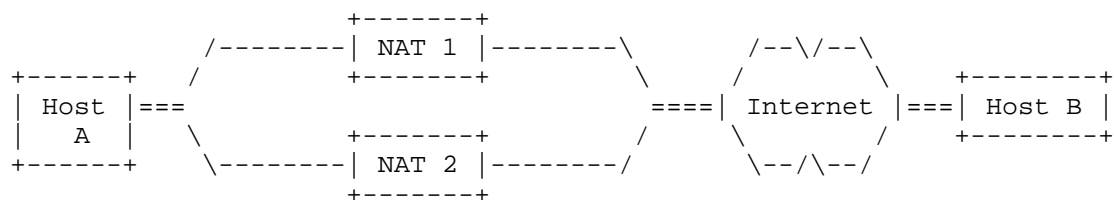
```

                                COOKIE-ACK
101.0.0.1:1 <----- 100.0.0.1:2
                        Int-VTag = 1234
```

```

                                COOKIE-ACK
10.0.0.1:1 <----- 100.0.0.1:2
      Int-VTag = 1234
```

Host A announces its second address in an ASCONF chunk. The address parameter contains an undefined address (0) to indicate that the source address should be added. The lookup address parameter within the ASCONF chunk will also contain the pair of VTags (external and internal) so that the NAT may populate its table completely with this single packet.



```

ASCONF [ADD-IP=0.0.0.0, INT-VTag=1234, Ext-VTag = 5678]
10.1.0.1:1 -----> 100.1.0.1:2
      Ext-VTag = 5678
```

NAT 2 creates complete entry:

NAT 2					
	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
	1234	1	10.1.0.1	5678	2

```

+-----+-----+-----+-----+-----+
ASCONF [ADD-IP,Int-VTag=1234, Ext-VTag = 5678]
101.1.0.1:1 -----> 100.1.0.1:2
                        Ext-VTag = 5678

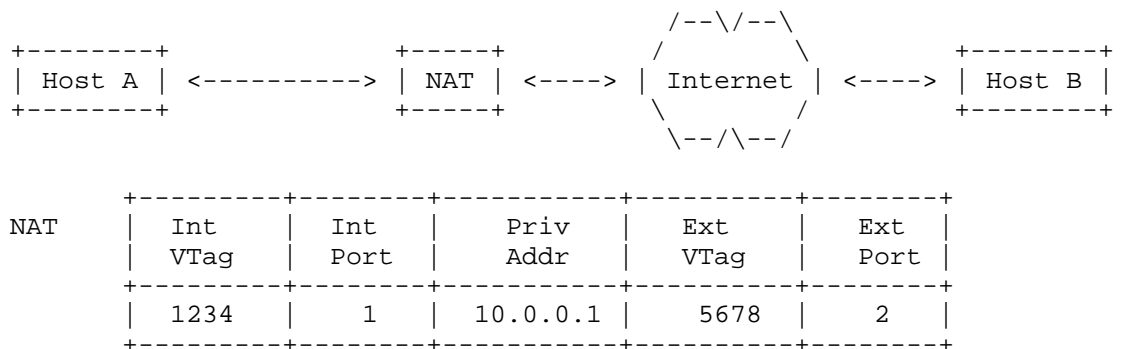
                        ASCONF-ACK
101.1.0.1:1 <----- 100.1.0.1:2
                        Int-VTag = 1234

ASCONF-ACK
10.1.0.1:1 <----- 100.1.0.1:2
      Int-VTag = 1234

```

9.4. NAT Loses Its State

Association is already established between Host A and Host B, when the NAT loses its state and obtains a new public address. Host A sends a DATA chunk to Host B.



```

      DATA
10.0.0.1:1 -----> 100.0.0.1:2
      Ext-VTag = 5678

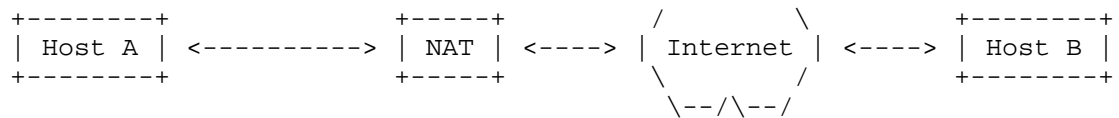
```

The NAT box cannot find entry for the association. It sends ERROR message with the M-Bit set and the cause "NAT state missing".

```

/--\ /--\

```

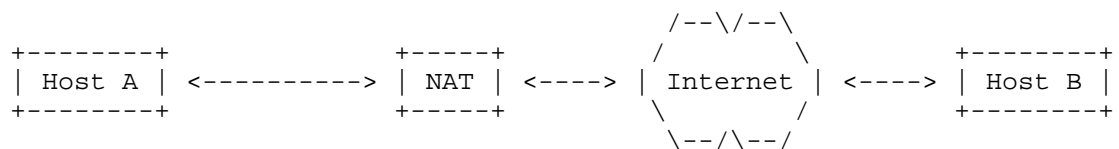


```

    ERROR [M-Bit, NAT state missing]
10.0.0.1:1 <----- 100.0.0.1:2
      Ext-VTag = 5678

```

On reception of the ERROR message, Host A sends an ASCONF chunk indicating that the former information has to be deleted and the source address of the actual packet added.



```

ASCONF [ADD-IP,DELETE-IP,Int-VTag=1234, Ext-VTag = 5678]
10.0.0.1:1 -----> 100.1.0.1:2
      Ext-VTag = 5678

```

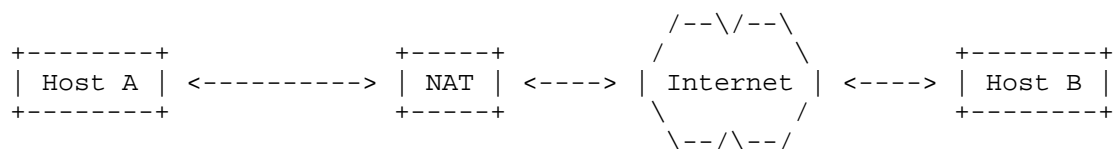
NAT	+-----+-----+-----+-----+-----+					
	Int	Int	Priv	Ext	Ext	
	VTag	Port	Addr	VTag	Port	
	+-----+-----+-----+-----+-----+					
	1234	1	10.0.0.1	5678	2	
	+-----+-----+-----+-----+-----+					

```

ASCONF [ADD-IP,DELETE-IP,Int-VTag=1234, Ext-VTag = 5678]
      102.1.0.1:1 -----> 100.1.0.1:2
                        Ext-VTag = 5678

```

Host B adds the new source address and deletes all former entries.




```

                                ASCONF-ACK
102.1.0.1:1 <----- 100.1.0.1:2
                                Int-VTag = 1234

                                ASCONF-ACK
10.1.0.1:1 <----- 100.1.0.1:2
                                Int-VTag = 1234

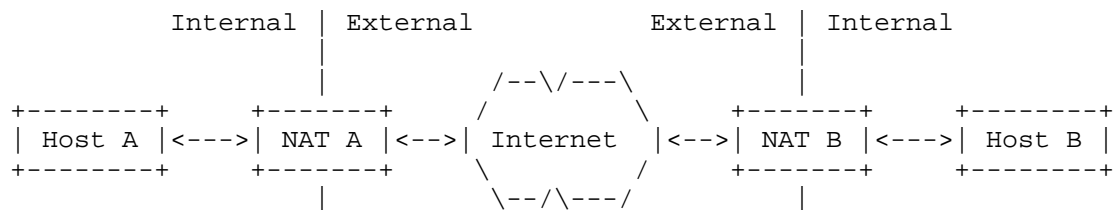
                                DATA
10.0.0.1:1 -----> 100.0.0.1:2
                                Ext-VTag = 5678

                                DATA
102.1.0.1:1 -----> 100.1.0.1:2
                                Ext-VTag = 5678

```

9.5. Peer-to-Peer Communication

If two hosts are behind NATs, they have to get knowledge of the peer's public address. This can be achieved with a so-called rendezvous server. Afterwards the destination addresses are public, and the association is set up with the help of the INIT collision. The NAT boxes create their entries according to their internal peer's point of view. Therefore, NAT A's Internal-VTag and Internal-Port are NAT B's External-VTag and External-Port, respectively. The naming of the verification tag in the packet flow is done from the sending peer's point of view.



NAT-Tables

NAT A	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
NAT B	Int v-tag	Int port	Priv addr	Ext v-tag	Ext port

```

+-----+-----+--- +-----+-----+
INIT[Initiate-Tag = 1234]
10.0.0.1:1 --> 100.0.0.1:2
      Ext-VTag = 0

```

NAT A creates entry:

NAT A					
	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
	1234	1	10.0.0.1	0	2

```

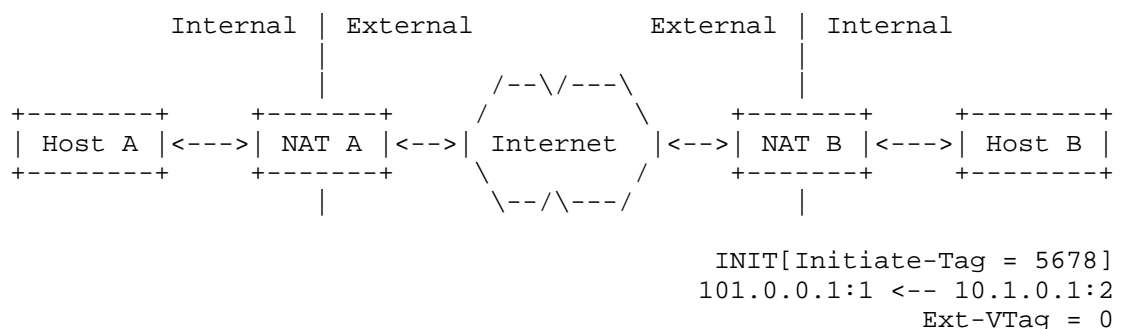
              INIT[Initiate-Tag = 1234]
101.0.0.1:1 -----> 100.0.0.1:2
              Ext-VTag = 0

```

NAT B processes INIT, but cannot find an entry. The SCTP packet is silently discarded and leaves the NAT table of NAT B unchanged.

NAT B					
	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port

Now Host B sends INIT, which is processed by NAT B. Its parameters are used to create an entry.



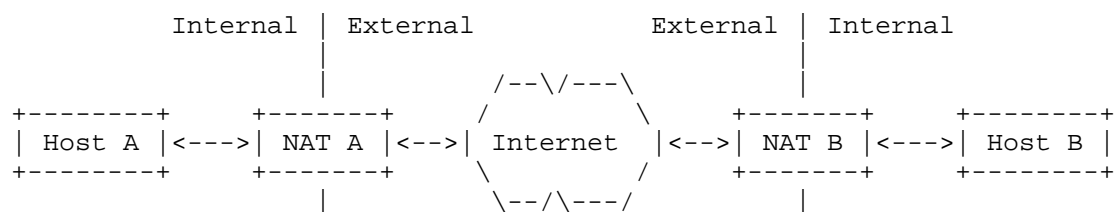
NAT B	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
	5678	2	10.1.0.1	0	1

```

INIT[Initiate-Tag = 5678]
101.0.0.1:1 <----- 100.0.0.1:2
                Ext-VTag = 0

```

NAT A processes INIT. As the outgoing INIT of Host A has already created an entry, the entry is found and updated:



VTag != Int-VTag, but Ext-VTag == 0, find entry.

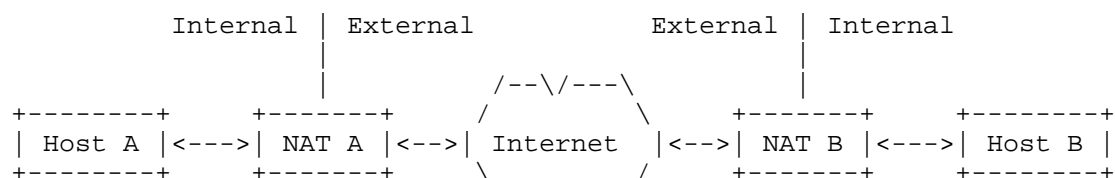
NAT A	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port
	1234	1	10.0.0.1	5678	2

```

INIT[Initiate-tag = 5678]
10.0.0.1:1 <-- 100.0.0.1:2
      Ext-VTag = 0

```

Host A send INIT-ACK, which can pass through NAT B:



```

|                               \--/\---/                               |
INIT-ACK[Initiate-Tag = 1234]
10.0.0.1:1 --> 100.0.0.1:2
    Ext-VTag = 5678

```

```

INIT-ACK[Initiate-Tag = 1234]
101.0.0.1:1 -----> 100.0.0.1:2
    Ext-VTag = 5678

```

NAT B updates entry:

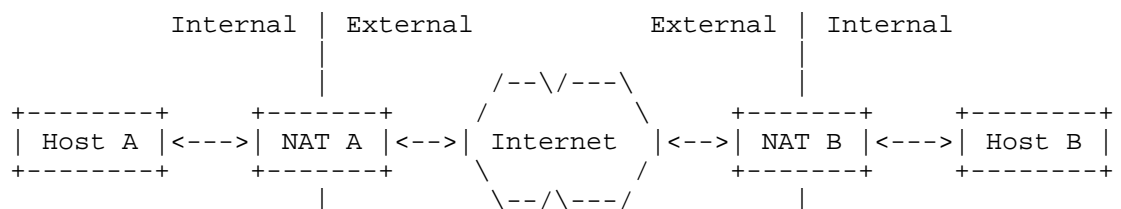
NAT B	+-----+-----+-----+-----+-----+					
	Int VTag	Int Port	Priv Addr	Ext VTag	Ext Port	
	5678	2	10.1.0.1	1234	1	

```

INIT-ACK[Initiate-Tag = 1234]
101.0.0.1:1 --> 10.1.0.1:2
    Ext-VTag = 5678

```

The lookup for COOKIE-ECHO and COOKIE-ACK is successful.



```

COOKIE-ECHO
101.0.0.1:1 <-- 10.1.0.1:2
    Ext-VTag = 1234

```

```

COOKIE-ECHO
101.0.0.1:1 <----- 100.0.0.1:2
    Ext-VTag = 1234

```

```

COOKIE-ECHO
10.0.0.1:1 <-- 100.0.0.1:2
    Ext-VTag = 1234

```

```
      COOKIE-ACK
10.0.0.1:1 --> 100.0.0.1:2
      Ext-VTag = 5678
```

```
      COOKIE-ACK
101.0.0.1:1 -----> 100.0.0.1:2
      Ext-VTag = 5678
```

```
      COOKIE-ACK
101.0.0.1:1 --> 10.1.0.1:2
      Ext-VTag = 5678
```

10. IANA Considerations

This document requires no actions from IANA.

11. Security Considerations

State maintenance within a NAT is always a subject of possible Denial Of Service attacks. This document recommends that at a minimum a NAT runs a timer on any SCTP state so that old association state can be cleaned up.

12. Acknowledgments

The authors wish to thank Jason But Bryan Ford, David Hayes, Alfred Hines, Henning Peters, Timo Voelker, Dan Wing, and Qiaobing Xie for their invaluable comments.

13. References

13.1. Normative References

- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, September 1981.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.
- [I-D.ietf-tsvwg-natsupp]

Stewart, R., Tuexen, M., and I. Ruengeler, "Stream Control Transmission Protocol (SCTP) Network Address Translation Support", draft-ietf-tsvwg-natsupp-05 (work in progress), February 2013.

13.2. Informative References

[RFC5735] Cotton, M. and L. Vegoda, "Special Use IPv4 Addresses", RFC 5735, January 2010.

Authors' Addresses

Randall R. Stewart
Adara Networks
Chapin, SC 29036
US

Email: randall@lakerest.net

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

Email: tuexen@fh-muenster.de

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

Email: i.ruengeler@fh-muenster.de

Transport Area Working Group
Internet-Draft
Updates: 2309 (if approved)
Intended status: BCP
Expires: May 11, 2014

B. Briscoe
BT
J. Manner
Aalto University
November 07, 2013

Byte and Packet Congestion Notification
draft-ietf-tsvwg-byte-pkt-congest-12

Abstract

This document provides recommendations of best current practice for dropping or marking packets using any active queue management (AQM) algorithm, including random early detection (RED), BLUE, pre-congestion notification (PCN) and newer schemes such as CoDel (Controlled Delay) and PIE (Proportional Integral controller Enhanced). We give three strong recommendations: (1) packet size should be taken into account when transports detect and respond to congestion indications, (2) packet size should not be taken into account when network equipment creates congestion signals (marking, dropping), and therefore (3) in the specific case of RED, the byte-mode packet drop variant that drops fewer small packets should not be used. This memo updates RFC 2309 to deprecate deliberate preferential treatment of small packets in AQM algorithms.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on May 11, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	4
1.1. Terminology and Scoping	6
1.2. Example Comparing Packet-Mode Drop and Byte-Mode Drop	7
2. Recommendations	9
2.1. Recommendation on Queue Measurement	9
2.2. Recommendation on Encoding Congestion Notification	10
2.3. Recommendation on Responding to Congestion	11
2.4. Recommendation on Handling Congestion Indications when Splitting or Merging Packets	12
3. Motivating Arguments	12
3.1. Avoiding Perverse Incentives to (Ab)use Smaller Packets	12
3.2. Small != Control	14
3.3. Transport-Independent Network	14
3.4. Partial Deployment of AQM	15
3.5. Implementation Efficiency	17
4. A Survey and Critique of Past Advice	17
4.1. Congestion Measurement Advice	18
4.1.1. Fixed Size Packet Buffers	18
4.1.2. Congestion Measurement without a Queue	19
4.2. Congestion Notification Advice	20
4.2.1. Network Bias when Encoding	20
4.2.2. Transport Bias when Decoding	22
4.2.3. Making Transports Robust against Control Packet Losses	23
4.2.4. Congestion Notification: Summary of Conflicting Advice	24
5. Outstanding Issues and Next Steps	25
5.1. Bit-congestible Network	25
5.2. Bit- & Packet-congestible Network	25
6. Security Considerations	26
7. IANA Considerations	26
8. Conclusions	26
9. Acknowledgements	28
10. Comments Solicited	28
11. References	28
11.1. Normative References	28
11.2. Informative References	28
Appendix A. Survey of RED Implementation Status	32
Appendix B. Sufficiency of Packet-Mode Drop	34
B.1. Packet-Size (In)Dependence in Transports	35
B.2. Bit-Congestible and Packet-Congestible Indications	38
Appendix C. Byte-mode Drop Complicates Policing Congestion Response	39
Appendix D. Changes from Previous Versions	40

1. Introduction

This document provides recommendations of best current practice for how we should correctly scale congestion control functions with respect to packet size for the long term. It also recognises that expediency may be necessary to deal with existing widely deployed protocols that don't live up to the long term goal.

When signalling congestion, the problem of how (and whether) to take packet sizes into account has exercised the minds of researchers and practitioners for as long as active queue management (AQM) has been discussed. Indeed, one reason AQM was originally introduced was to reduce the lock-out effects that small packets can have on large packets in drop-tail queues. This memo aims to state the principles we should be using and to outline how these principles will affect future protocol design, taking into account the existing deployments we have already.

The question of whether to take into account packet size arises at three stages in the congestion notification process:

Measuring congestion: When a congested resource measures locally how congested it is, should it measure its queue length in time, bytes or packets?

Encoding congestion notification into the wire protocol: When a congested network resource signals its level of congestion, should it drop / mark each packet dependent on the size of the particular packet in question?

Decoding congestion notification from the wire protocol: When a transport interprets the notification in order to decide how much to respond to congestion, should it take into account the size of each missing or marked packet?

Consensus has emerged over the years concerning the first stage, which Section 2.1 records in the RFC Series. In summary: If possible it is best to measure congestion by time in the queue, but otherwise the choice between bytes and packets solely depends on whether the resource is congested by bytes or packets.

The controversy is mainly around the last two stages: whether to allow for the size of the specific packet notifying congestion i) when the network encodes or ii) when the transport decodes the congestion notification.

Currently, the RFC series is silent on this matter other than a paper trail of advice referenced from [RFC2309], which conditionally

recommends byte-mode (packet-size dependent) drop [pktByteEmail]. Reducing drop of small packets certainly has some tempting advantages: i) it drops less control packets, which tend to be small and ii) it makes TCP's bit-rate less dependent on packet size. However, there are ways of addressing these issues at the transport layer, rather than reverse engineering network forwarding to fix the problems.

This memo updates [RFC2309] to deprecate deliberate preferential treatment of packets in AQM algorithms solely because of their size. It recommends that (1) packet size should be taken into account when transports detect and respond to congestion indications, (2) not when network equipment creates them. This memo also adds to the congestion control principles enumerated in BCP 41 [RFC2914].

In the particular case of Random early Detection (RED), this means that the byte-mode packet drop variant should not be used to drop fewer small packets, because that creates a perverse incentive for transports to use tiny segments, consequently also opening up a DoS vulnerability. Fortunately all the RED implementers who responded to our admittedly limited survey (Section 4.2.4) have not followed the earlier advice to use byte-mode drop, so the position this memo argues for seems to already exist in implementations.

However, at the transport layer, TCP congestion control is a widely deployed protocol that doesn't scale with packet size (i.e. its reduction in rate does not take into account the size of a lost packet). To date this hasn't been a significant problem because most TCP implementations have been used with similar packet sizes. But, as we design new congestion control mechanisms, this memo recommends that we should build in scaling with packet size rather than assuming we should follow TCP's example.

This memo continues as follows. First it discusses terminology and scoping. Section 2 gives the concrete formal recommendations, followed by motivating arguments in Section 3. We then critically survey the advice given previously in the RFC series and the research literature (Section 4), referring to an assessment of whether or not this advice has been followed in production networks (Appendix A). To wrap up, outstanding issues are discussed that will need resolution both to inform future protocol designs and to handle legacy (Section 5). Then security issues are collected together in Section 6 before conclusions are drawn in Section 8. The interested reader can find discussion of more detailed issues on the theme of byte vs. packet in the appendices.

This memo intentionally includes a non-negligible amount of material on the subject. For the busy reader Section 2 summarises the

recommendations for the Internet community.

1.1. Terminology and Scoping

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].

This memo applies to the design of all AQM algorithms, for example, Random Early Detection (RED) [RFC2309], BLUE [BLUE02], Pre-Congestion Notification (PCN) [RFC5670], Controlled Delay (CoDel) [I-D.nichols-tsvwg-codel] and the Proportional Integral controller Enhanced (PIE) [I-D.pan-tsvwg-pie]. Throughout, RED is used as a concrete example because it is a widely known and deployed AQM algorithm. There is no intention to imply that the advice is any less applicable to the other algorithms, nor that RED is preferred.

Congestion Notification: Congestion notification is a changing signal that aims to communicate the probability that the network resource(s) will not be able to forward the level of traffic load offered (or that there is an impending risk that they will not be able to).

The 'impending risk' qualifier is added, because AQM systems set a virtual limit smaller than the actual limit to the resource, then notify when this virtual limit is exceeded in order to avoid uncontrolled congestion of the actual capacity.

Congestion notification communicates a real number bounded by the range [0 , 1]. This ties in with the most well-understood measure of congestion notification: drop probability.

Explicit and Implicit Notification: The byte vs. packet dilemma concerns congestion notification irrespective of whether it is signalled implicitly by drop or using Explicit Congestion Notification (ECN [RFC3168] or PCN [RFC5670]). Throughout this document, unless clear from the context, the term marking will be used to mean notifying congestion explicitly, while congestion notification will be used to mean notifying congestion either implicitly by drop or explicitly by marking.

Bit-congestible vs. Packet-congestible: If the load on a resource depends on the rate at which packets arrive, it is called packet-congestible. If the load depends on the rate at which bits arrive it is called bit-congestible.

Examples of packet-congestible resources are route look-up engines and firewalls, because load depends on how many packet headers

they have to process. Examples of bit-congestible resources are transmission links, radio power and most buffer memory, because the load depends on how many bits they have to transmit or store. Some machine architectures use fixed size packet buffers, so buffer memory in these cases is packet-congestible (see Section 4.1.1).

The path through a machine will typically encounter both packet-congestible and bit-congestible resources. However, currently, a design goal of network processing equipment such as routers and firewalls is to size the packet-processing engine(s) relative to the lines in order to keep packet processing uncongested even under worst case packet rates with runs of minimum size packets. Therefore, packet-congestion is currently rare [RFC6077; S.3.3], but there is no guarantee that it will not become more common in future.

Note that information is generally processed or transmitted with a minimum granularity greater than a bit (e.g. octets). The appropriate granularity for the resource in question should be used, but for the sake of brevity we will talk in terms of bytes in this memo.

Coarser Granularity: Resources may be congestible at higher levels of granularity than bits or packets, for instance stateful firewalls are flow-congestible and call-servers are session-congestible. This memo focuses on congestion of connectionless resources, but the same principles may be applicable for congestion notification protocols controlling per-flow and per-session processing or state.

RED Terminology: In RED whether to use packets or bytes when measuring queues is called respectively "packet-mode queue measurement" or "byte-mode queue measurement". And whether the probability of dropping a particular packet is independent or dependent on its size is called respectively "packet-mode drop" or "byte-mode drop". The terms byte-mode and packet-mode should not be used without specifying whether they apply to queue measurement or to drop.

1.2. Example Comparing Packet-Mode Drop and Byte-Mode Drop

Taking RED as a well-known example algorithm, a central question addressed by this document is whether to recommend RED's packet-mode drop variant and to deprecate byte-mode drop. Table 1 compares how packet-mode and byte-mode drop affect two flows of different size packets. For each it gives the expected number of packets and of bits dropped in one second. Each example flow runs at the same bit-

rate of 48Mb/s, but one is broken up into small 60 byte packets and the other into large 1500 byte packets.

To keep up the same bit-rate, in one second there are about 25 times more small packets because they are 25 times smaller. As can be seen from the table, the packet rate is 100,000 small packets versus 4,000 large packets per second (pps).

Parameter	Formula	Small packets	Large packets
Packet size	$s/8$	60B	1,500B
Packet size	s	480b	12,000b
Bit-rate	x	48Mbps	48Mbps
Packet-rate	$u = x/s$	100kpps	4kpps
Packet-mode Drop			
Pkt loss probability	p	0.1%	0.1%
Pkt loss-rate	$p*u$	100pps	4pps
Bit loss-rate	$p*u*s$	48kbps	48kbps
Byte-mode Drop			
	MTU, $M=12,000b$		
Pkt loss probability	$b = p*s/M$	0.004%	0.1%
Pkt loss-rate	$b*u$	4pps	4pps
Bit loss-rate	$b*u*s$	1.92kbps	48kbps

Table 1: Example Comparing Packet-mode and Byte-mode Drop

For packet-mode drop, we illustrate the effect of a drop probability of 0.1%, which the algorithm applies to all packets irrespective of size. Because there are 25 times more small packets in one second, it naturally drops 25 times more small packets, that is 100 small packets but only 4 large packets. But if we count how many bits it drops, there are 48,000 bits in 100 small packets and 48,000 bits in 4 large packets--the same number of bits of small packets as large.

The packet-mode drop algorithm drops any bit with the same probability whether the bit is in a small or a large packet.

For byte-mode drop, again we use an example drop probability of 0.1%, but only for maximum size packets (assuming the link maximum transmission unit (MTU) is 1,500B or 12,000b). The byte-mode algorithm reduces the drop probability of smaller packets proportional to their size, making the probability that it drops a small packet 25 times smaller at 0.004%. But there are 25 times more small packets, so dropping them with 25 times lower probability results in dropping the same number of packets: 4 drops in both cases. The 4 small dropped packets contain 25 times less bits than the 4 large dropped packets: 1,920 compared to 48,000.

The byte-mode drop algorithm drops any bit with a probability proportionate to the size of the packet it is in.

2. Recommendations

This section gives recommendations related to network equipment in Sections 2.1 and 2.2, and in Sections 2.3 and 2.4 we discuss the implications on the transport protocols.

2.1. Recommendation on Queue Measurement

Ideally, an AQM would measure the service time of the queue to measure congestion of a resource. However service time can only be measured as packets leave the queue, where it is not always expedient to implement a full AQM algorithm. To predict the service time as packets join the queue, an AQM algorithm needs to measure the length of the queue.

In this case, if the resource is bit-congestible, the AQM implementation SHOULD measure the length of the queue in bytes and, if the resource is packet-congestible, the implementation SHOULD measure the length of the queue in packets. Subject to the exceptions below, no other choice makes sense, because the number of packets waiting in the queue isn't relevant if the resource gets congested by bytes and vice versa. For example, the length of the queue into a transmission line would be measured in bytes, while the length of the queue into a firewall would be measured in packets.

To avoid the pathological effects of drop tail, the AQM can then transform this service time or queue length into the probability of dropping or marking a packet (e.g. RED's piecewise linear function between thresholds).

What this advice means for RED as a specific example:

1. A RED implementation SHOULD use byte mode queue measurement for measuring the congestion of bit-congestible resources and packet mode queue measurement for packet-congestible resources.
2. An implementation SHOULD NOT make it possible to configure the way a queue measures itself, because whether a queue is bit-congestible or packet-congestible is an inherent property of the queue.

Exceptions to these recommendations might be necessary, for instance where a packet-congestible resource has to be configured as a proxy bottleneck for a bit-congestible resource in an adjacent box that does not support AQM.

The recommended approach in less straightforward scenarios, such as fixed size packet buffers, resources without a queue and buffers comprising a mix of packet and bit-congestible resources, is discussed in Section 4.1. For instance, Section 4.1.1 explains that the queue into a line should be measured in bytes even if the queue consists of fixed-size packet-buffers, because the root-cause of any congestion is bytes arriving too fast for the line--packets filling buffers are merely a symptom of the underlying congestion of the line.

2.2. Recommendation on Encoding Congestion Notification

When encoding congestion notification (e.g. by drop, ECN or PCN), the probability that network equipment drops or marks a particular packet to notify congestion SHOULD NOT depend on the size of the packet in question. As the example in Section 1.2 illustrates, to drop any bit with probability 0.1% it is only necessary to drop every packet with probability 0.1% without regard to the size of each packet.

This approach ensures the network layer offers sufficient congestion information for all known and future transport protocols and also ensures no perverse incentives are created that would encourage transports to use inappropriately small packet sizes.

What this advice means for RED as a specific example:

1. The RED AQM algorithm SHOULD NOT use byte-mode drop, i.e. it ought to use packet-mode drop. Byte-mode drop is more complex, it creates the perverse incentive to fragment segments into tiny pieces and it is vulnerable to floods of small packets.
2. If a vendor has implemented byte-mode drop, and an operator has turned it on, it is RECOMMENDED to switch it to packet-mode drop, after establishing if there are any implications on the relative performance of applications using different packet sizes. The unlikely possibility of some application-specific legacy use of byte-mode drop is the only reason that all the above recommendations on encoding congestion notification are not phrased more strongly.

RED as a whole SHOULD NOT be switched off. Without RED, a drop tail queue biases against large packets and is vulnerable to floods of small packets.

Note well that RED's byte-mode queue drop is completely orthogonal to byte-mode queue measurement and should not be confused with it. If a RED implementation has a byte-mode but does not specify what sort of byte-mode, it is most probably byte-mode queue measurement, which is

fine. However, if in doubt, the vendor should be consulted.

A survey (Appendix A) showed that there appears to be little, if any, installed base of the byte-mode drop variant of RED. This suggests that deprecating byte-mode drop will have little, if any, incremental deployment impact.

2.3. Recommendation on Responding to Congestion

When a transport detects that a packet has been lost or congestion marked, it SHOULD consider the strength of the congestion indication as proportionate to the size in octets (bytes) of the missing or marked packet.

In other words, when a packet indicates congestion (by being lost or marked) it can be considered conceptually as if there is a congestion indication on every octet of the packet, not just one indication per packet.

To be clear, the above recommendation solely describes how a transport should interpret the meaning of a congestion indication, as a long term goal. It makes no recommendation on whether a transport should act differently based on this interpretation. It merely aids interoperability between transports, if they choose to make their actions depend on the strength of congestion indications.

This definition will be useful as the IETF transport area continues its programme of;

- o updating host-based congestion control protocols to take account of packet size
- o making transports less sensitive to losing control packets like SYNs and pure ACKs.

What this advice means for the case of TCP:

1. If two TCP flows with different packet sizes are required to run at equal bit rates under the same path conditions, this SHOULD be done by altering TCP (Section 4.2.2), not network equipment (the latter affects other transports besides TCP).
2. If it is desired to improve TCP performance by reducing the chance that a SYN or a pure ACK will be dropped, this SHOULD be done by modifying TCP (Section 4.2.3), not network equipment.

To be clear, we are not recommending at all that TCPs under equivalent conditions should aim for equal bit-rates. We are merely

saying that anyone trying to do such a thing should modify their TCP algorithm, not the network.

These recommendations are phrased as 'SHOULD' rather than 'MUST', because there may be cases where expediency dictates that compatibility with pre-existing versions of a transport protocol make the recommendations impractical.

2.4. Recommendation on Handling Congestion Indications when Splitting or Merging Packets

Packets carrying congestion indications may be split or merged in some circumstances (e.g. at a RTP/RTCP transcoder or during IP fragment reassembly). Splitting and merging only make sense in the context of ECN, not loss.

The general rule to follow is that the number of octets in packets with congestion indications SHOULD be equivalent before and after merging or splitting. This is based on the principle used above; that an indication of congestion on a packet can be considered as an indication of congestion on each octet of the packet.

The above rule is not phrased with the word "MUST" to allow the following exception. There are cases where pre-existing protocols were not designed to conserve congestion marked octets (e.g. IP fragment reassembly [RFC3168] or loss statistics in RTCP receiver reports [RFC3550] before ECN was added [RFC6679]). When any such protocol is updated, it SHOULD comply with the above rule to conserve marked octets. However, the rule may be relaxed if it would otherwise become too complex to interoperate with pre-existing implementations of the protocol.

One can think of a splitting or merging process as if all the incoming congestion-marked octets increment a counter and all the outgoing marked octets decrement the same counter. In order to ensure that congestion indications remain timely, even the smallest positive remainder in the conceptual counter should trigger the next outgoing packet to be marked (causing the counter to go negative).

3. Motivating Arguments

This section is informative. It justifies the recommendations given in the previous section.

3.1. Avoiding Perverse Incentives to (Ab)use Smaller Packets

Increasingly, it is being recognised that a protocol design must take care not to cause unintended consequences by giving the parties in

the protocol exchange perverse incentives [Evol_cc][RFC3426]. Given there are many good reasons why larger path maximum transmission units (PMTUs) would help solve a number of scaling issues, we do not want to create any bias against large packets that is greater than their true cost.

Imagine a scenario where the same bit rate of packets will contribute the same to bit-congestion of a link irrespective of whether it is sent as fewer larger packets or more smaller packets. A protocol design that caused larger packets to be more likely to be dropped than smaller ones would be dangerous in both the following cases:

Malicious transports: A queue that gives an advantage to small packets can be used to amplify the force of a flooding attack. By sending a flood of small packets, the attacker can get the queue to discard more traffic in large packets, allowing more attack traffic to get through to cause further damage. Such a queue allows attack traffic to have a disproportionately large effect on regular traffic without the attacker having to do much work.

Non-malicious transports: Even if an application designer is not actually malicious, if over time it is noticed that small packets tend to go faster, designers will act in their own interest and use smaller packets. Queues that give advantage to small packets create an evolutionary pressure for applications or transports to send at the same bit-rate but break their data stream down into tiny segments to reduce their drop rate. Encouraging a high volume of tiny packets might in turn unnecessarily overload a completely unrelated part of the system, perhaps more limited by header-processing than bandwidth.

Imagine two unresponsive flows arrive at a bit-congestible transmission link each with the same bit rate, say 1Mbps, but one consists of 1500B and the other 60B packets, which are 25x smaller. Consider a scenario where gentle RED [gentle_RED] is used, along with the variant of RED we advise against, i.e. where the RED algorithm is configured to adjust the drop probability of packets in proportion to each packet's size (byte mode packet drop). In this case, RED aims to drop 25x more of the larger packets than the smaller ones. Thus, for example if RED drops 25% of the larger packets, it will aim to drop 1% of the smaller packets (but in practice it may drop more as congestion increases [RFC4828; Appx B.4]). Even though both flows arrive with the same bit rate, the bit rate the RED queue aims to pass to the line will be 750kbps for the flow of larger packets but 990kbps for the smaller packets (because of rate variations it will actually be a little less than this target).

Note that, although the byte-mode drop variant of RED amplifies small

packet attacks, drop-tail queues amplify small packet attacks even more (see Security Considerations in Section 6). Wherever possible neither should be used.

3.2. Small != Control

Dropping fewer control packets considerably improves performance. It is tempting to drop small packets with lower probability in order to improve performance, because many control packets tend to be smaller (TCP SYNs & ACKs, DNS queries & responses, SIP messages, HTTP GETs, etc). However, we must not give control packets preference purely by virtue of their smallness, otherwise it is too easy for any data source to get the same preferential treatment simply by sending data in smaller packets. Again we should not create perverse incentives to favour small packets rather than to favour control packets, which is what we intend.

Just because many control packets are small does not mean all small packets are control packets.

So, rather than fix these problems in the network, we argue that the transport should be made more robust against losses of control packets (see 'Making Transports Robust against Control Packet Losses' in Section 4.2.3).

3.3. Transport-Independent Network

TCP congestion control ensures that flows competing for the same resource each maintain the same number of segments in flight, irrespective of segment size. So under similar conditions, flows with different segment sizes will get different bit-rates.

To counter this effect it seems tempting not to follow our recommendation, and instead for the network to bias congestion notification by packet size in order to equalise the bit-rates of flows with different packet sizes. However, in order to do this, the queuing algorithm has to make assumptions about the transport, which become embedded in the network. Specifically:

- o The queuing algorithm has to assume how aggressively the transport will respond to congestion (see Section 4.2.4). If the network assumes the transport responds as aggressively as TCP NewReno, it will be wrong for Compound TCP and differently wrong for Cubic TCP, etc. To achieve equal bit-rates, each transport then has to guess what assumption the network made, and work out how to replace this assumed aggressiveness with its own aggressiveness.

- o Also, if the network biases congestion notification by packet size it has to assume a baseline packet size--all proposed algorithms use the local MTU (for example see the byte-mode loss probability formula in Table 1). Then if the non-Reno transports mentioned above are trying to reverse engineer what the network assumed, they also have to guess the MTU of the congested link.

Even though reducing the drop probability of small packets (e.g. RED's byte-mode drop) helps ensure TCP flows with different packet sizes will achieve similar bit rates, we argue this correction should be made to any future transport protocols based on TCP, not to the network in order to fix one transport, no matter how predominant it is. Effectively, favouring small packets is reverse engineering of network equipment around one particular transport protocol (TCP), contrary to the excellent advice in [RFC3426], which asks designers to question "Why are you proposing a solution at this layer of the protocol stack, rather than at another layer?"

In contrast, if the network never takes account of packet size, the transport can be certain it will never need to guess any assumptions the network has made. And the network passes two pieces of information to the transport that are sufficient in all cases: i) congestion notification on the packet and ii) the size of the packet. Both are available for the transport to combine (by taking account of packet size when responding to congestion) or not. Appendix B checks that these two pieces of information are sufficient for all relevant scenarios.

When the network does not take account of packet size, it allows transport protocols to choose whether to take account of packet size or not. However, if the network were to bias congestion notification by packet size, transport protocols would have no choice; those that did not take account of packet size themselves would unwittingly become dependent on packet size, and those that already took account of packet size would end up taking account of it twice.

3.4. Partial Deployment of AQM

In overview, the argument in this section runs as follows:

- o Because the network does not and cannot always drop packets in proportion to their size, it shouldn't be given the task of making drop signals depend on packet size at all.
- o Transports on the other hand don't always want to make their rate response proportional to the size of dropped packets, but if they want to, they always can.

The argument is similar to the end-to-end argument that says "Don't do X in the network if end-systems can do X by themselves, and they want to be able to choose whether to do X anyway." Actually the following argument is stronger; in addition it says "Don't give the network task X that could be done by the end-systems, if X is not deployed on all network nodes, and end-systems won't be able to tell whether their network is doing X, or whether they need to do X themselves." In this case, the X in question is "making the response to congestion depend on packet size".

We will now re-run this argument taking each step in more depth. The argument applies solely to drop, not to ECN marking.

A queue drops packets for either of two reasons: a) to signal to host congestion controls that they should reduce the load and b) because there is no buffer left to store the packets. Active queue management tries to use drops as a signal for hosts to slow down (case a) so that drop due to buffer exhaustion (case b) should not be necessary.

AQM is not universally deployed in every queue in the Internet; many cheap Ethernet bridges, software firewalls, NATs on consumer devices, etc implement simple tail-drop buffers. Even if AQM were universal, it has to be able to cope with buffer exhaustion (by switching to a behaviour like tail-drop), in order to cope with unresponsive or excessive transports. For these reasons networks will sometimes be dropping packets as a last resort (case b) rather than under AQM control (case a).

When buffers are exhausted (case b), they don't naturally drop packets in proportion to their size. The network can only reduce the probability of dropping smaller packets if it has enough space to store them somewhere while it waits for a larger packet that it can drop. If the buffer is exhausted, it does not have this choice. Admittedly tail-drop does naturally drop somewhat fewer small packets, but exactly how few depends more on the mix of sizes than the size of the packet in question. Nonetheless, in general, if we wanted networks to do size-dependent drop, we would need universal deployment of (packet-size dependent) AQM code, which is currently unrealistic.

A host transport cannot know whether any particular drop was a deliberate signal from an AQM or a sign of a queue shedding packets due to buffer exhaustion. Therefore, because the network cannot universally do size-dependent drop, it should not do it all.

Whereas universality is desirable in the network, diversity is desirable between different transport layer protocols - some, like

NewReno TCP [RFC5681], may not choose to make their rate response proportionate to the size of each dropped packet, while others will (e.g. TFRC-SP [RFC4828]).

3.5. Implementation Efficiency

Biasing against large packets typically requires an extra multiply and divide in the network (see the example byte-mode drop formula in Table 1). Allowing for packet size at the transport rather than in the network ensures that neither the network nor the transport needs to do a multiply operation--multiplication by packet size is effectively achieved as a repeated add when the transport adds to its count of marked bytes as each congestion event is fed to it. Also the work to do the biasing is spread over many hosts, rather than concentrated in just the congested network element. These aren't principled reasons in themselves, but they are a happy consequence of the other principled reasons.

4. A Survey and Critique of Past Advice

This section is informative, not normative.

The original 1993 paper on RED [RED93] proposed two options for the RED active queue management algorithm: packet mode and byte mode. Packet mode measured the queue length in packets and dropped (or marked) individual packets with a probability independent of their size. Byte mode measured the queue length in bytes and marked an individual packet with probability in proportion to its size (relative to the maximum packet size). In the paper's outline of further work, it was stated that no recommendation had been made on whether the queue size should be measured in bytes or packets, but noted that the difference could be significant.

When RED was recommended for general deployment in 1998 [RFC2309], the two modes were mentioned implying the choice between them was a question of performance, referring to a 1997 email [pktByteEmail] for advice on tuning. A later addendum to this email introduced the insight that there are in fact two orthogonal choices:

- o whether to measure queue length in bytes or packets (Section 4.1)
- o whether the drop probability of an individual packet should depend on its own size (Section 4.2).

The rest of this section is structured accordingly.

4.1. Congestion Measurement Advice

The choice of which metric to use to measure queue length was left open in RFC2309. It is now well understood that queues for bit-congestible resources should be measured in bytes, and queues for packet-congestible resources should be measured in packets [pktByteEmail].

Congestion in some legacy bit-congestible buffers is only measured in packets not bytes. In such cases, the operator has to set the thresholds mindful of a typical mix of packets sizes. Any AQM algorithm on such a buffer will be oversensitive to high proportions of small packets, e.g. a DoS attack, and under-sensitive to high proportions of large packets. However, there is no need to make allowances for the possibility of such legacy in future protocol design. This is safe because any under-sensitivity during unusual traffic mixes cannot lead to congestion collapse given the buffer will eventually revert to tail drop, discarding proportionately more large packets.

4.1.1. Fixed Size Packet Buffers

The question of whether to measure queues in bytes or packets seems to be well understood. However, measuring congestion is confusing when the resource is bit congestible but the queue into the resource is packet congestible. This section outlines the approach to take.

Some, mostly older, queuing hardware allocates fixed sized buffers in which to store each packet in the queue. This hardware forwards to the line in one of two ways:

- o With some hardware, any fixed sized buffers not completely filled by a packet are padded when transmitted to the wire. This case, should clearly be treated as packet-congestible, because both queuing and transmission are in fixed MTU-sized units. Therefore the queue length in packets is a good model of congestion of the link.
- o More commonly, hardware with fixed size packet buffers transmits packets to line without padding. This implies a hybrid forwarding system with transmission congestion dependent on the size of packets but queue congestion dependent on the number of packets, irrespective of their size.

Nonetheless, there would be no queue at all unless the line had become congested--the root-cause of any congestion is too many bytes arriving for the line. Therefore, the AQM should measure the queue length as the sum of all the packet sizes in bytes that

are queued up waiting to be serviced by the line, irrespective of whether each packet is held in a fixed size buffer.

In the (unlikely) first case where use of padding means the queue should be measured in packets, further confusion is likely because the fixed buffers are rarely all one size. Typically pools of different sized buffers are provided (Cisco uses the term 'buffer carving' for the process of dividing up memory into these pools [IOSArch]). Usually, if the pool of small buffers is exhausted, arriving small packets can borrow space in the pool of large buffers, but not vice versa. However, there is no need to consider all this complexity, because the root-cause of any congestion is still line overload--buffer consumption is only the symptom. Therefore, the length of the queue should be measured as the sum of the bytes in the queue that will be transmitted to line, including any padding. In the (unusual) case of transmission with padding this means the sum of the sizes of the small buffers queued plus the sum of the sizes of the large buffers queued.

We will return to borrowing of fixed sized buffers when we discuss biasing the drop/marketing probability of a specific packet because of its size in Section 4.2.1. But here we can repeat the simple rule for how to measure the length of queues of fixed buffers: no matter how complicated the buffering scheme is, ultimately a transmission line is nearly always bit-congestible so the number of bytes queued up waiting for the line measures how congested the line is, and it is rarely important to measure how congested the buffering system is.

4.1.1.2. Congestion Measurement without a Queue

AQM algorithms are nearly always described assuming there is a queue for a congested resource and the algorithm can use the queue length to determine the probability that it will drop or mark each packet. But not all congested resources lead to queues. For instance, power limited resources are usually bit-congestible if energy is primarily required for transmission rather than header processing, but it is rare for a link protocol to build a queue as it approaches maximum power.

Nonetheless, AQM algorithms do not require a queue in order to work. For instance spectrum congestion can be modelled by signal quality using target bit-energy-to-noise-density ratio. And, to model radio power exhaustion, transmission power levels can be measured and compared to the maximum power available. [ECNFixedWireless] proposes a practical and theoretically sound way to combine congestion notification for different bit-congestible resources at different layers along an end to end path, whether wireless or wired, and whether with or without queues.

In wireless protocols that use request to send / clear to send (RTS / CTS) control, such as some variants of IEEE802.11, it is reasonable to base an AQM on the time spent waiting for transmission opportunities (TXOPs) even though wireless spectrum is usually regarded as congested by bits (for a given coding scheme). This is because requests for TXOPs queue up as the spectrum gets congested by all the bits being transferred. So the time that TXOPs are queued directly reflects bit congestion of the spectrum.

4.2. Congestion Notification Advice

4.2.1. Network Bias when Encoding

4.2.1.1. Advice on Packet Size Bias in RED

The previously mentioned email [pktByteEmail] referred to by [RFC2309] advised that most scarce resources in the Internet were bit-congestible, which is still believed to be true (Section 1.1). But it went on to offer advice that is updated by this memo. It said that drop probability should depend on the size of the packet being considered for drop if the resource is bit-congestible, but not if it is packet-congestible. The argument continued that if packet drops were inflated by packet size (byte-mode dropping), "a flow's fraction of the packet drops is then a good indication of that flow's fraction of the link bandwidth in bits per second". This was consistent with a referenced policing mechanism being worked on at the time for detecting unusually high bandwidth flows, eventually published in 1999 [pBox]. However, the problem could and should have been solved by making the policing mechanism count the volume of bytes randomly dropped, not the number of packets.

A few months before RFC2309 was published, an addendum was added to the above archived email referenced from the RFC, in which the final paragraph seemed to partially retract what had previously been said. It clarified that the question of whether the probability of dropping/markings a packet should depend on its size was not related to whether the resource itself was bit congestible, but a completely orthogonal question. However the only example given had the queue measured in packets but packet drop depended on the size of the packet in question. No example was given the other way round.

In 2000, Cnodder et al [REDbyte] pointed out that there was an error in the part of the original 1993 RED algorithm that aimed to distribute drops uniformly, because it didn't correctly take into account the adjustment for packet size. They recommended an algorithm called RED_4 to fix this. But they also recommended a further change, RED_5, to adjust drop rate dependent on the square of relative packet size. This was indeed consistent with one implied

motivation behind RED's byte mode drop--that we should reverse engineer the network to improve the performance of dominant end-to-end congestion control mechanisms. This memo makes a different recommendations in Section 2.

By 2003, a further change had been made to the adjustment for packet size, this time in the RED algorithm of the ns2 simulator. Instead of taking each packet's size relative to a 'maximum packet size' it was taken relative to a 'mean packet size', intended to be a static value representative of the 'typical' packet size on the link. We have not been able to find a justification in the literature for this change, however Eddy and Allman conducted experiments [REDbias] that assessed how sensitive RED was to this parameter, amongst other things. However, this changed algorithm can often lead to drop probabilities of greater than 1 (which gives a hint that there is probably a mistake in the theory somewhere).

On 10-Nov-2004, this variant of byte-mode packet drop was made the default in the ns2 simulator. It seems unlikely that byte-mode drop has ever been implemented in production networks (Appendix A), therefore any conclusions based on ns2 simulations that use RED without disabling byte-mode drop are likely to behave very differently from RED in production networks.

4.2.1.2. Packet Size Bias Regardless of AQM

The byte-mode drop variant of RED (or a similar variant of other AQM algorithms) is not the only possible bias towards small packets in queueing systems. We have already mentioned that tail-drop queues naturally tend to lock-out large packets once they are full.

But also queues with fixed sized buffers reduce the probability that small packets will be dropped if (and only if) they allow small packets to borrow buffers from the pools for larger packets (see Section 4.1.1). Borrowing effectively makes the maximum queue size for small packets greater than that for large packets, because more buffers can be used by small packets while less will fit large packets. Incidentally, the bias towards small packets from buffer borrowing is nothing like as large as that of RED's byte-mode drop.

Nonetheless, fixed-buffer memory with tail drop is still prone to lock-out large packets, purely because of the tail-drop aspect. So, fixed size packet-buffers should be augmented with a good AQM algorithm and packet-mode drop. If an AQM is too complicated to implement with multiple fixed buffer pools, the minimum necessary to prevent large packet lock-out is to ensure smaller packets never use the last available buffer in any of the pools for larger packets.

4.2.2. Transport Bias when Decoding

The above proposals to alter the network equipment to bias towards smaller packets have largely carried on outside the IETF process. Whereas, within the IETF, there are many different proposals to alter transport protocols to achieve the same goals, i.e. either to make the flow bit-rate take account of packet size, or to protect control packets from loss. This memo argues that altering transport protocols is the more principled approach.

A recently approved experimental RFC adapts its transport layer protocol to take account of packet sizes relative to typical TCP packet sizes. This proposes a new small-packet variant of TCP-friendly rate control [RFC5348] called TFRC-SP [RFC4828]. Essentially, it proposes a rate equation that inflates the flow rate by the ratio of a typical TCP segment size (1500B including TCP header) over the actual segment size [PktSizeEquCC]. (There are also other important differences of detail relative to TFRC, such as using virtual packets [CCvarPktSize] to avoid responding to multiple losses per round trip and using a minimum inter-packet interval.)

Section 4.5.1 of this TFRC-SP spec discusses the implications of operating in an environment where queues have been configured to drop smaller packets with proportionately lower probability than larger ones. But it only discusses TCP operating in such an environment, only mentioning TFRC-SP briefly when discussing how to define fairness with TCP. And it only discusses the byte-mode dropping version of RED as it was before Cnoddler et al pointed out it didn't sufficiently bias towards small packets to make TCP independent of packet size.

So the TFRC-SP spec doesn't address the issue of which of the network or the transport should handle fairness between different packet sizes. In its Appendix B.4 it discusses the possibility of both TFRC-SP and some network buffers duplicating each other's attempts to deliberately bias towards small packets. But the discussion is not conclusive, instead reporting simulations of many of the possibilities in order to assess performance but not recommending any particular course of action.

The paper originally proposing TFRC with virtual packets (VP-TFRC) [CCvarPktSize] proposed that there should perhaps be two variants to cater for the different variants of RED. However, as the TFRC-SP authors point out, there is no way for a transport to know whether some queues on its path have deployed RED with byte-mode packet drop (except if an exhaustive survey found that no-one has deployed it!--see Appendix A). Incidentally, VP-TFRC also proposed that byte-mode RED dropping should really square the packet-size compensation-factor

(like that of Cnoder's RED_5, but apparently unaware of it).

Pre-congestion notification [RFC5670] is an IETF technology to use a virtual queue for AQM marking for packets within one Diffserv class in order to give early warning prior to any real queuing. The PCN marking algorithms have been designed not to take account of packet size when forwarding through queues. Instead the general principle has been to take account of the sizes of marked packets when monitoring the fraction of marking at the edge of the network, as recommended here.

4.2.3. Making Transports Robust against Control Packet Losses

Recently, two RFCs have defined changes to TCP that make it more robust against losing small control packets [RFC5562] [RFC5690]. In both cases they note that the case for these two TCP changes would be weaker if RED were biased against dropping small packets. We argue here that these two proposals are a safer and more principled way to achieve TCP performance improvements than reverse engineering RED to benefit TCP.

Although there are no known proposals, it would also be possible and perfectly valid to make control packets robust against drop by requesting a scheduling class with lower drop probability, by re-marking to a Diffserv code point [RFC2474] within the same behaviour aggregate.

Although not brought to the IETF, a simple proposal from Wischik [DupTCP] suggests that the first three packets of every TCP flow should be routinely duplicated after a short delay. It shows that this would greatly improve the chances of short flows completing quickly, but it would hardly increase traffic levels on the Internet, because Internet bytes have always been concentrated in the large flows. It further shows that the performance of many typical applications depends on completion of long serial chains of short messages. It argues that, given most of the value people get from the Internet is concentrated within short flows, this simple expedient would greatly increase the value of the best efforts Internet at minimal cost. A similar but more extensive approach has been evaluated on Google servers [GentleAggro].

The proposals discussed in this sub-section are experimental approaches that are not yet in wide operational use, but they are existence proofs that transports can make themselves robust against loss of control packets. The examples are all TCP-based, but applications over non-TCP transports could mitigate loss of control packets by making similar use of Diffserv, data duplication, FEC etc.

4.2.4. Congestion Notification: Summary of Conflicting Advice

transport cc	RED_1 (packet mode drop)	RED_4 (linear byte mode drop)	RED_5 (square byte mode drop)
TCP or TFRC	s/\sqrt{p}	$\sqrt{s/p}$	$1/\sqrt{p}$
TFRC-SP	$1/\sqrt{p}$	$1/\sqrt{sp}$	$1/(s.\sqrt{p})$

Table 2: Dependence of flow bit-rate per RTT on packet size, s , and drop probability, p , when network and/or transport bias towards small packets to varying degrees

Table 2 aims to summarise the potential effects of all the advice from different sources. Each column shows a different possible AQM behaviour in different queues in the network, using the terminology of Cnoder et al outlined earlier (RED_1 is basic RED with packet-mode drop). Each row shows a different transport behaviour: TCP [RFC5681] and TFRC [RFC5348] on the top row with TFRC-SP [RFC4828] below. Each cell shows how the bits per round trip of a flow depends on packet size, s , and drop probability, p . In order to declutter the formulae to focus on packet-size dependence they are all given per round trip, which removes any RTT term.

Let us assume that the goal is for the bit-rate of a flow to be independent of packet size. Suppressing all inessential details, the table shows that this should either be achievable by not altering the TCP transport in a RED_5 network, or using the small packet TFRC-SP transport (or similar) in a network without any byte-mode dropping RED (top right and bottom left). Top left is the 'do nothing' scenario, while bottom right is the 'do-both' scenario in which bit-rate would become far too biased towards small packets. Of course, if any form of byte-mode dropping RED has been deployed on a subset of queues that congest, each path through the network will present a different hybrid scenario to its transport.

Whatever, we can see that the linear byte-mode drop column in the middle would considerably complicate the Internet. It's a half-way house that doesn't bias enough towards small packets even if one believes the network should be doing the biasing. Section 2 recommends that all bias in network equipment towards small packets should be turned off--if indeed any equipment vendors have implemented it--leaving packet-size bias solely as the preserve of the transport layer (solely the leftmost, packet-mode drop column).

In practice it seems that no deliberate bias towards small packets

has been implemented for production networks. Of the 19% of vendors who responded to a survey of 84 equipment vendors, none had implemented byte-mode drop in RED (see Appendix A for details).

5. Outstanding Issues and Next Steps

5.1. Bit-congestible Network

For a connectionless network with nearly all resources being bit-congestible the recommended position is clear--that the network should not make allowance for packet sizes and the transport should. This leaves two outstanding issues:

- o How to handle any legacy of AQM with byte-mode drop already deployed;
- o The need to start a programme to update transport congestion control protocol standards to take account of packet size.

A survey of equipment vendors (Section 4.2.4) found no evidence that byte-mode packet drop had been implemented, so deployment will be sparse at best. A migration strategy is not really needed to remove an algorithm that may not even be deployed.

A programme of experimental updates to take account of packet size in transport congestion control protocols has already started with TFRC-SP [RFC4828].

5.2. Bit- & Packet-congestible Network

The position is much less clear-cut if the Internet becomes populated by a more even mix of both packet-congestible and bit-congestible resources (see Appendix B.2). This problem is not pressing, because most Internet resources are designed to be bit-congestible before packet processing starts to congest (see Section 1.1).

The IRTF Internet congestion control research group (ICCRG) has set itself the task of reaching consensus on generic forwarding mechanisms that are necessary and sufficient to support the Internet's future congestion control requirements (the first challenge in [RFC6077]). The research question of whether packet congestion might become common and what to do if it does may in the future be explored in the IRTF (the "Challenge 3: Packet Size" in [RFC6077]).

Note that sometimes it seems that resources might be congested by neither bits nor packets, e.g. where the queue for access to a wireless medium is in units of transmission opportunities. However,

the root cause of congestion of the underlying spectrum is overload of bits (see Section 4.1.2).

6. Security Considerations

This memo recommends that queues do not bias drop probability due to packets size. For instance dropping small packets less often than large creates a perverse incentive for transports to break down their flows into tiny segments. One of the benefits of implementing AQM was meant to be to remove this perverse incentive that drop-tail queues gave to small packets.

In practice, transports cannot all be trusted to respond to congestion. So another reason for recommending that queues do not bias drop probability towards small packets is to avoid the vulnerability to small packet DDoS attacks that would otherwise result. One of the benefits of implementing AQM was meant to be to remove drop-tail's DoS vulnerability to small packets, so we shouldn't add it back again.

If most queues implemented AQM with byte-mode drop, the resulting network would amplify the potency of a small packet DDoS attack. At the first queue the stream of packets would push aside a greater proportion of large packets, so more of the small packets would survive to attack the next queue. Thus a flood of small packets would continue on towards the destination, pushing regular traffic with large packets out of the way in one queue after the next, but suffering much less drop itself.

Appendix C explains why the ability of networks to police the response of any transport to congestion depends on bit-congestible network resources only doing packet-mode not byte-mode drop. In summary, it says that making drop probability depend on the size of the packets that bits happen to be divided into simply encourages the bits to be divided into smaller packets. Byte-mode drop would therefore irreversibly complicate any attempt to fix the Internet's incentive structures.

7. IANA Considerations

This document has no actions for IANA.

8. Conclusions

This memo identifies the three distinct stages of the congestion notification process where implementations need to decide whether to take packet size into account. The recommendations provided in Section 2 of this memo are different in each case:

- o When network equipment measures the length of a queue, if it is not feasible to use time it is recommended to count in bytes if the network resource is congested by bytes, or to count in packets if is congested by packets.
- o When network equipment decides whether to drop (or mark) a packet, it is recommended that the size of the particular packet should not be taken into account
- o However, when a transport algorithm responds to a dropped or marked packet, the size of the rate reduction should be proportionate to the size of the packet.

In summary, the answers are 'it depends', 'no' and 'yes' respectively

For the specific case of RED, this means that byte-mode queue measurement will often be appropriate but the use of byte-mode drop is very strongly discouraged.

At the transport layer the IETF should continue updating congestion control protocols to take account of the size of each packet that indicates congestion. Also the IETF should continue to make protocols less sensitive to losing control packets like SYN's, pure ACKs and DNS exchanges. Although many control packets happen to be small, the alternative of network equipment favouring all small packets would be dangerous. That would create perverse incentives to split data transfers into smaller packets.

The memo develops these recommendations from principled arguments concerning scaling, layering, incentives, inherent efficiency, security and policeability. But it also addresses practical issues such as specific buffer architectures and incremental deployment. Indeed a limited survey of RED implementations is discussed, which shows there appears to be little, if any, installed base of RED's byte-mode drop. Therefore it can be deprecated with little, if any, incremental deployment complications.

The recommendations have been developed on the well-founded basis that most Internet resources are bit-congestible not packet-congestible. We need to know the likelihood that this assumption will prevail longer term and, if it might not, what protocol changes will be needed to cater for a mix of the two. The IRTF Internet Congestion Control Research Group (ICCRG) is currently working on these problems [RFC6077].

9. Acknowledgements

Thank you to Sally Floyd, who gave extensive and useful review comments. Also thanks for the reviews from Philip Eardley, David Black, Fred Baker, David Taht, Toby Moncaster, Arnaud Jacquet and Mirja Kuehlewind as well as helpful explanations of different hardware approaches from Larry Dunn and Fred Baker. We are grateful to Bruce Davie and his colleagues for providing a timely and efficient survey of RED implementation in Cisco's product range. Also grateful thanks to Toby Moncaster, Will Dormann, John Regnault, Simon Carter and Stefaan De Cnodder who further helped survey the current status of RED implementation and deployment and, finally, thanks to the anonymous individuals who responded.

Bob Briscoe and Jukka Manner were partly funded by Trilogy, a research project (ICT- 216372) supported by the European Community under its Seventh Framework Programme. The views expressed here are those of the authors only.

10. Comments Solicited

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

11. References

11.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.

11.2. Informative References

- [BLUE02] Feng, W-c., Shin, K., Kandlur, D., and D. Saha, "The BLUE active queue management algorithms", IEEE/ACM Transactions on Networking 10(4) 513--528, August 2002, <http://dx.doi.org/10.1109/TNET.2002.801399>.
- [CCvarPktSize] Widmer, J., Boutremans, C., and J-Y. Le

- Boudec, "Congestion Control for Flows with Variable Packet Size", ACM CCR 34(2) 137--151, 2004, <<http://doi.acm.org/10.1145/997150.997162>>.
- [CHOKE_Var_Pkt] Psounis, K., Pan, R., and B. Prabhaker, "Approximate Fair Dropping for Variable Length Packets", IEEE Micro 21(1):48--56, January-February 2001, <<http://www.stanford.edu/~balaji/papers/01approximatefair.pdf>>.
- [DRQ] Shin, M., Chong, S., and I. Rhee, "Dual-Resource TCP/AQM for Processing-Constrained Networks", IEEE/ACM Transactions on Networking Vol 16, issue 2, April 2008, <<http://dx.doi.org/10.1109/TNET.2007.900415>>.
- [DupTCP] Wischik, D., "Short messages", Philosophical Transactions of the Royal Society A 366(1872):1941-1953, June 2008, <<http://rsta.royalsocietypublishing.org/content/366/1872/1941.full.pdf+html>>.
- [ECNFixedWireless] Siris, V., "Resource Control for Elastic Traffic in CDMA Networks", Proc. ACM MOBICOM'02 , September 2002, <http://www.ics.forth.gr/netlab/publications/resource_control_elastic_cdma.html>.
- [Evol_cc] Gibbens, R. and F. Kelly, "Resource pricing and the evolution of congestion control", Automatica 35(12):1969--1985, December 1999, <<http://www.statslab.cam.ac.uk/~frank/evol.html>>.
- [GentleAggro] Flach, T., Dukkupati, N., Terzis, A., Raghavan, B., Cardwell, N., Cheng, Y., Jain, A., Hao, S., Katz-Bassett, E., and R. Govindan, "Reducing Web Latency: the Virtue of Gentle Aggression", ACM SIGCOMM CCR 43(4):159--170, August 2013, <<http://doi.acm.org/10.1145/2486001.2486014>>.
- [I-D.nichols-tsvwg-codel] Nichols, K. and V. Jacobson, "Controlled Delay Active Queue Management",

- draft-nichols-tsvwg-codel-01 (work in progress), February 2013.
- [I-D.pan-tsvwg-pie] Pan, R., Natarajan, P., Piglione, C., and M. Prabhu, "PIE: A Lightweight Control Scheme To Address the Bufferbloat Problem", draft-pan-tsvwg-pie-00 (work in progress), December 2012.
- [IOSArch] Bollapragada, V., White, R., and C. Murphy, "Inside Cisco IOS Software Architecture", Cisco Press: CCIE Professional Development ISBN13: 978-1-57870-181-0, July 2000.
- [PktSizeEquCC] Vasallo, P., "Variable Packet Size Equation-Based Congestion Control", ICSI Technical Report tr-00-008, 2000, <<http://http.icsi.berkeley.edu/ftp/global/pub/techreports/2000/tr-00-008.pdf>>.
- [RED93] Floyd, S. and V. Jacobson, "Random Early Detection (RED) gateways for Congestion Avoidance", IEEE/ACM Transactions on Networking 1(4) 397--413, August 1993, <<http://www.icir.org/floyd/papers/red/red.html>>.
- [REDBias] Eddy, W. and M. Allman, "A Comparison of RED's Byte and Packet Modes", Computer Networks 42(3) 261--280, June 2003, <<http://www.ir.bbn.com/documents/articles/redbias.ps>>.
- [REDbyte] De Cnodder, S., Elloumi, O., and K. Pauwels, "RED behavior with different packet sizes", Proc. 5th IEEE Symposium on Computers and Communications (ISCC) 793--799, July 2000, <<http://www.icir.org/floyd/red/Elloumi99.pdf>>.
- [RFC2309] Braden, B., Clark, D., Crowcroft, J., Davie, B., Deering, S., Estrin, D., Floyd, S., Jacobson, V., Minshall, G., Partridge, C., Peterson, L., Ramakrishnan, K., Shenker, S., Wroclawski, J., and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet",

RFC 2309, April 1998.

- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, December 1998.
- [RFC2914] Floyd, S., "Congestion Control Principles", BCP 41, RFC 2914, September 2000.
- [RFC3426] Floyd, S., "General Architectural and Policy Considerations", RFC 3426, November 2002.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3714] Floyd, S. and J. Kempf, "IAB Concerns Regarding Congestion Control for Voice Traffic in the Internet", RFC 3714, March 2004.
- [RFC4828] Floyd, S. and E. Kohler, "TCP Friendly Rate Control (TFRC): The Small-Packet (SP) Variant", RFC 4828, April 2007.
- [RFC5348] Floyd, S., Handley, M., Padhye, J., and J. Widmer, "TCP Friendly Rate Control (TFRC): Protocol Specification", RFC 5348, September 2008.
- [RFC5562] Kuzmanovic, A., Mondal, A., Floyd, S., and K. Ramakrishnan, "Adding Explicit Congestion Notification (ECN) Capability to TCP's SYN/ACK Packets", RFC 5562, June 2009.
- [RFC5670] Eardley, P., "Metering and Marking Behaviour of PCN-Nodes", RFC 5670, November 2009.
- [RFC5681] Allman, M., Paxson, V., and E. Blanton, "TCP Congestion Control", RFC 5681, September 2009.

- [RFC5690] Floyd, S., Arcia, A., Ros, D., and J. Iyengar, "Adding Acknowledgement Congestion Control to TCP", RFC 5690, February 2010.
- [RFC6077] Papadimitriou, D., Welzl, M., Scharf, M., and B. Briscoe, "Open Research Issues in Internet Congestion Control", RFC 6077, February 2011.
- [RFC6679] Westerlund, M., Johansson, I., Perkins, C., O'Hanlon, P., and K. Carlberg, "Explicit Congestion Notification (ECN) for RTP over UDP", RFC 6679, August 2012.
- [RFC6789] Briscoe, B., Woundy, R., and A. Cooper, "Congestion Exposure (ConEx) Concepts and Use Cases", RFC 6789, December 2012.
- [Rate_fair_Dis] Briscoe, B., "Flow Rate Fairness: Dismantling a Religion", ACM CCR 37(2)63--74, April 2007, <<http://portal.acm.org/citation.cfm?id=1232926>>.
- [gentle_RED] Floyd, S., "Recommendation on using the "gentle_" variant of RED", Web page , March 2000, <<http://www.icir.org/floyd/red/gentle.html>>.
- [pBox] Floyd, S. and K. Fall, "Promoting the Use of End-to-End Congestion Control in the Internet", IEEE/ACM Transactions on Networking 7(4) 458--472, August 1999, <<http://www.aciri.org/floyd/end2end-paper.html>>.
- [pktByteEmail] Floyd, S., "RED: Discussions of Byte and Packet Modes", email , March 1997, <<http://www-nrg.ee.lbl.gov/floyd/REDaveraging.txt>>.

Appendix A. Survey of RED Implementation Status

This Appendix is informative, not normative.

In May 2007 a survey was conducted of 84 vendors to assess how widely drop probability based on packet size has been implemented in RED Table 3. About 19% of those surveyed replied, giving a sample size

of 16. Although in most cases we do not have permission to identify the respondents, we can say that those that have responded include most of the larger equipment vendors, covering a large fraction of the market. The two who gave permission to be identified were Cisco and Alcatel-Lucent. The others range across the large network equipment vendors at L3 & L2, firewall vendors, wireless equipment vendors, as well as large software businesses with a small selection of networking products. All those who responded confirmed that they have not implemented the variant of RED with drop dependent on packet size (2 were fairly sure they had not but needed to check more thoroughly). At the time the survey was conducted, Linux did not implement RED with packet-size bias of drop, although we have not investigated a wider range of open source code.

Response	No. of vendors	%age of vendors
Not implemented	14	17%
Not implemented (probably)	2	2%
Implemented	0	0%
No response	68	81%
Total companies/orgs surveyed	84	100%

Table 3: Vendor Survey on byte-mode drop variant of RED (lower drop probability for small packets)

Where reasons have been given, the extra complexity of packet bias code has been most prevalent, though one vendor had a more principled reason for avoiding it--similar to the argument of this document.

Our survey was of vendor implementations, so we cannot be certain about operator deployment. But we believe many queues in the Internet are still tail-drop. The company of one of the co-authors (BT) has widely deployed RED, but many tail-drop queues are bound to still exist, particularly in access network equipment and on middleboxes like firewalls, where RED is not always available.

Routers using a memory architecture based on fixed size buffers with borrowing may also still be prevalent in the Internet. As explained in Section 4.2.1, these also provide a marginal (but legitimate) bias towards small packets. So even though RED byte-mode drop is not prevalent, it is likely there is still some bias towards small packets in the Internet due to tail drop and fixed buffer borrowing.

Appendix B. Sufficiency of Packet-Mode Drop

This Appendix is informative, not normative.

Here we check that packet-mode drop (or marking) in the network gives sufficiently generic information for the transport layer to use. We check against a 2x2 matrix of four scenarios that may occur now or in the future (Table 4). The horizontal and vertical dimensions have been chosen because each tests extremes of sensitivity to packet size in the transport and in the network respectively.

Note that this section does not consider byte-mode drop at all. Having deprecated byte-mode drop, the goal here is to check that packet-mode drop will be sufficient in all cases.

Network	Transport	a) Independent of packet size of congestion notifications	b) Dependent on packet size of congestion notifications
1) Predominantly bit-congestible network		Scenario a1)	Scenario b1)
2) Mix of bit-congestible and pkt-congestible network		Scenario a2)	Scenario b2)

Table 4: Four Possible Congestion Scenarios

Appendix B.1 focuses on the horizontal dimension of Table 4 checking that packet-mode drop (or marking) gives sufficient information, whether or not the transport uses it--scenarios b) and a) respectively.

Appendix B.2 focuses on the vertical dimension of Table 4, checking that packet-mode drop gives sufficient information to the transport whether resources in the network are bit-congestible or packet-congestible (these terms are defined in Section 1.1).

Notation: To be concrete, we will compare two flows with different packet sizes, s_1 and s_2 . As an example, we will take $s_1 = 60B = 480b$ and $s_2 = 1500B = 12,000b$.

A flow's bit rate, x [bps], is related to its packet rate, u [pps], by

$$x(t) = s.u(t).$$

In the bit-congestible case, path congestion will be denoted by `p_b`, and in the packet-congestible case by `p_p`. When either case is implied, the letter `p` alone will denote path congestion.

B.1. Packet-Size (In)Dependence in Transports

In all cases we consider a packet-mode drop queue that indicates congestion by dropping (or marking) packets with probability `p` irrespective of packet size. We use an example value of loss (marking) probability, `p=0.1%`.

A transport like RFC5681 TCP treats a congestion notification on any packet whatever its size as one event. However, a network with just the packet-mode drop algorithm does give more information if the transport chooses to use it. We will use Table 5 to illustrate this.

We will set aside the last column until later. The columns labelled "Flow 1" and "Flow 2" compare two flows consisting of 60B and 1500B packets respectively. The body of the table considers two separate cases, one where the flows have equal bit-rate and the other with equal packet-rates. In both cases, the two flows fill a 96Mbps link. Therefore, in the equal bit-rate case they each have half the bit-rate (48Mbps). Whereas, with equal packet-rates, flow 1 uses 25 times smaller packets so it gets 25 times less bit-rate--it only gets $1/(1+25)$ of the link capacity ($96\text{Mbps}/26 = 4\text{Mbps}$ after rounding). In contrast flow 2 gets 25 times more bit-rate (92Mbps) in the equal packet rate case because its packets are 25 times larger. The packet rate shown for each flow could easily be derived once the bit-rate was known by dividing bit-rate by packet size, as shown in the column labelled "Formula".

Parameter	Formula	Flow 1	Flow 2	Combined
-----	-----	-----	-----	-----
Packet size	$s/8$	60B	1,500B	(Mix)
Packet size	s	480b	12,000b	(Mix)
Pkt loss probability	p	0.1%	0.1%	0.1%
EQUAL BIT-RATE CASE				
Bit-rate	x	48Mbps	48Mbps	96Mbps
Packet-rate	$u = x/s$	100kpps	4kpps	104kpps
Absolute pkt-loss-rate	$p*u$	100pps	4pps	104pps
Absolute bit-loss-rate	$p*u*s$	48kbps	48kbps	96kbps
Ratio of lost/sent pkts	$p*u/u$	0.1%	0.1%	0.1%
Ratio of lost/sent bits	$p*u*s/(u*s)$	0.1%	0.1%	0.1%
EQUAL PACKET-RATE CASE				
Bit-rate	x	4Mbps	92Mbps	96Mbps
Packet-rate	$u = x/s$	8kpps	8kpps	15kpps
Absolute pkt-loss-rate	$p*u$	8pps	8pps	15pps
Absolute bit-loss-rate	$p*u*s$	4kbps	92kbps	96kbps
Ratio of lost/sent pkts	$p*u/u$	0.1%	0.1%	0.1%
Ratio of lost/sent bits	$p*u*s/(u*s)$	0.1%	0.1%	0.1%

Table 5: Absolute Loss Rates and Loss Ratios for Flows of Small and Large Packets and Both Combined

So far we have merely set up the scenarios. We now consider congestion notification in the scenario. Two TCP flows with the same round trip time aim to equalise their packet-loss-rates over time. That is the number of packets lost in a second, which is the packets per second (u) multiplied by the probability that each one is dropped (p). Thus TCP converges on the "Equal packet-rate" case, where both flows aim for the same "Absolute packet-loss-rate" (both 8pps in the table).

Packet-mode drop actually gives flows sufficient information to measure their loss-rate in bits per second, if they choose, not just packets per second. Each flow can count the size of a lost or marked packet and scale its rate-response in proportion (as TFRC-SP does). The result is shown in the row entitled "Absolute bit-loss-rate", where the bits lost in a second is the packets per second (u) multiplied by the probability of losing a packet (p) multiplied by the packet size (s). Such an algorithm would try to remove any imbalance in bit-loss-rate such as the wide disparity in the "Equal packet-rate" case (4kbps vs. 92kbps). Instead, a packet-size-dependent algorithm would aim for equal bit-loss-rates, which would drive both flows towards the "Equal bit-rate" case, by driving them to equal bit-loss-rates (both 48kbps in this example).

The explanation so far has assumed that each flow consists of packets of only one constant size. Nonetheless, it extends naturally to flows with mixed packet sizes. In the right-most column of Table 5 a flow of mixed size packets is created simply by considering flow 1 and flow 2 as a single aggregated flow. There is no need for a flow to maintain an average packet size. It is only necessary for the transport to scale its response to each congestion indication by the size of each individual lost (or marked) packet. Taking for example the "Equal packet-rate" case, in one second about 8 small packets and 8 large packets are lost (making closer to 15 than 16 losses per second due to rounding). If the transport multiplies each loss by its size, in one second it responds to $8 \times 480\text{b}$ and $8 \times 12,000\text{b}$ lost bits, adding up to 96,000 lost bits in a second. This double checks correctly, being the same as 0.1% of the total bit-rate of 96Mbps. For completeness, the formula for absolute bit-loss-rate is $p(u_1 \times s_1 + u_2 \times s_2)$.

Incidentally, a transport will always measure the loss probability the same irrespective of whether it measures in packets or in bytes. In other words, the ratio of lost to sent packets will be the same as the ratio of lost to sent bytes. (This is why TCP's bit rate is still proportional to packet size even when byte-counting is used, as recommended for TCP in [RFC5681], mainly for orthogonal security reasons.) This is intuitively obvious by comparing two example flows; one with 60B packets, the other with 1500B packets. If both flows pass through a queue with drop probability 0.1%, each flow will lose 1 in 1,000 packets. In the stream of 60B packets the ratio of bytes lost to sent will be 60B in every 60,000B; and in the stream of 1500B packets, the loss ratio will be 1,500B out of 1,500,000B. When the transport responds to the ratio of lost to sent packets, it will measure the same ratio whether it measures in packets or bytes: 0.1% in both cases. The fact that this ratio is the same whether measured in packets or bytes can be seen in Table 5, where the ratio of lost to sent packets and the ratio of lost to sent bytes is always 0.1% in all cases (recall that the scenario was set up with $p=0.1\%$).

This discussion of how the ratio can be measured in packets or bytes is only raised here to highlight that it is irrelevant to this memo! Whether a transport depends on packet size or not depends on how this ratio is used within the congestion control algorithm.

So far we have shown that packet-mode drop passes sufficient information to the transport layer so that the transport can take account of bit-congestion, by using the sizes of the packets that indicate congestion. We have also shown that the transport can choose not to take packet size into account if it wishes. We will now consider whether the transport can know which to do.

B.2. Bit-Congestible and Packet-Congestible Indications

As a thought-experiment, imagine an idealised congestion notification protocol that supports both bit-congestible and packet-congestible resources. It would require at least two ECN flags, one for each of bit-congestible and packet-congestible resources.

1. A packet-congestible resource trying to code congestion level p_p into a packet stream should mark the idealised 'packet congestion' field in each packet with probability p_p irrespective of the packet's size. The transport should then take a packet with the packet congestion field marked to mean just one mark, irrespective of the packet size.
2. A bit-congestible resource trying to code time-varying byte-congestion level p_b into a packet stream should mark the 'byte congestion' field in each packet with probability p_b , again irrespective of the packet's size. Unlike before, the transport should take a packet with the byte congestion field marked to count as a mark on each byte in the packet.

This hides a fundamental problem--much more fundamental than whether we can magically create header space for yet another ECN flag, or whether it would work while being deployed incrementally. Distinguishing drop from delivery naturally provides just one implicit bit of congestion indication information--the packet is either dropped or not. It is hard to drop a packet in two ways that are distinguishable remotely. This is a similar problem to that of distinguishing wireless transmission losses from congestive losses.

This problem would not be solved even if ECN were universally deployed. A congestion notification protocol must survive a transition from low levels of congestion to high. Marking two states is feasible with explicit marking, but much harder if packets are dropped. Also, it will not always be cost-effective to implement AQM at every low level resource, so drop will often have to suffice.

We are not saying two ECN fields will be needed (and we are not saying that somehow a resource should be able to drop a packet in one of two different ways so that the transport can distinguish which sort of drop it was!). These two congestion notification channels are a conceptual device to illustrate a dilemma we could face in the future. Section 3 gives four good reasons why it would be a bad idea to allow for packet size by biasing drop probability in favour of small packets within the network. The impracticality of our thought experiment shows that it will be hard to give transports a practical way to know whether to take account of the size of congestion indication packets or not.

Fortunately, this dilemma is not pressing because by design most equipment becomes bit-congested before its packet-processing becomes congested (as already outlined in Section 1.1). Therefore transports can be designed on the relatively sound assumption that a congestion indication will usually imply bit-congestion.

Nonetheless, although the above idealised protocol isn't intended for implementation, we do want to emphasise that research is needed to predict whether there are good reasons to believe that packet congestion might become more common, and if so, to find a way to somehow distinguish between bit and packet congestion [RFC3714].

Recently, the dual resource queue (DRQ) proposal [DRQ] has been made on the premise that, as network processors become more cost effective, per packet operations will become more complex (irrespective of whether more function in the network is desirable). Consequently the premise is that CPU congestion will become more common. DRQ is a proposed modification to the RED algorithm that folds both bit congestion and packet congestion into one signal (either loss or ECN).

Finally, we note one further complication. Strictly, packet-congestible resources are often cycle-congestible. For instance, for routing look-ups load depends on the complexity of each look-up and whether the pattern of arrivals is amenable to caching or not. This also reminds us that any solution must not require a forwarding engine to use excessive processor cycles in order to decide how to say it has no spare processor cycles.

Appendix C. Byte-mode Drop Complicates Policing Congestion Response

This section is informative, not normative.

There are two main classes of approach to policing congestion response: i) policing at each bottleneck link or ii) policing at the edges of networks. Packet-mode drop in RED is compatible with either, while byte-mode drop precludes edge policing.

The simplicity of an edge policer relies on one dropped or marked packet being equivalent to another of the same size without having to know which link the drop or mark occurred at. However, the byte-mode drop algorithm has to depend on the local MTU of the line--it needs to use some concept of a 'normal' packet size. Therefore, one dropped or marked packet from a byte-mode drop algorithm is not necessarily equivalent to another from a different link. A policing function local to the link can know the local MTU where the congestion occurred. However, a policer at the edge of the network cannot, at least not without a lot of complexity.

The early research proposals for type (i) policing at a bottleneck link [pBox] used byte-mode drop, then detected flows that contributed disproportionately to the number of packets dropped. However, with no extra complexity, later proposals used packet mode drop and looked for flows that contributed a disproportionate amount of dropped bytes [CHOKe_Var_Pkt].

Work is progressing on the congestion exposure protocol (ConEx [RFC6789]), which enables a type (ii) edge policer located at a user's attachment point. The idea is to be able to take an integrated view of the effect of all a user's traffic on any link in the internetwork. However, byte-mode drop would effectively preclude such edge policing because of the MTU issue above.

Indeed, making drop probability depend on the size of the packets that bits happen to be divided into would simply encourage the bits to be divided into smaller packets in order to confuse policing. In contrast, as long as a dropped/marked packet is taken to mean that all the bytes in the packet are dropped/marked, a policer can remain robust against bits being re-divided into different size packets or across different size flows [Rate_fair_Dis].

Appendix D. Changes from Previous Versions

To be removed by the RFC Editor on publication.

Full incremental diffs between each version are available at
<<http://tools.ietf.org/wg/tsvwg/draft-ietf-tsvwg-byte-pkt-congest/>>
(courtesy of the rfcdiff tool):

From -11 to -12: Following the second pass through the IESG:

- * Section 2.1 [Barry Leiba]:
 - + s/No other choice makes sense,/Subject to the exceptions below, no other choice makes sense,/
 - + s/Exceptions to these recommendations MAY be necessary /Exceptions to these recommendations may be necessary /
- * Sections 3.2 and 4.2.3 [Joel Jaeggli]:
 - + Added comment to section 4.2.3 that the examples given are not in widespread production use, but they give evidence that it is possible to follow the advice given.
 - + Section 4.2.3:

- OLD: Although there are no known proposals, it would also be possible and perfectly valid to make control packets robust against drop by explicitly requesting a lower drop probability using their Diffserv code point [RFC2474] to request a scheduling class with lower drop.
NEW: Although there are no known proposals, it would also be possible and perfectly valid to make control packets robust against drop by requesting a scheduling class with lower drop probability, by re-marking to a Diffserv code point [RFC2474] within the same behaviour aggregate.
- appended "Similarly applications, over non-TCP transports could make any packets that are effectively control packets more robust by using Diffserv, data duplication, FEC etc."
- + Updated Wischik ref and added "Reducing Web Latency: the Virtue of Gentle Aggression" ref.
- * Expanded more abbreviations (CoDel, PIE, MTU).
- * Section 1. Intro [Stephen Farrell]:
 - + In the places where the doc describes the dichotomy between 'long-term goal' and 'expediency' the words long term goal and expedient have been introduced, to more explicitly refer back to this introductory para (S.2.1 & S.2.3).
 - + Added explanation of what scaling with packet size means.
- * Conclusions [Benoit Claise]:
 - + OLD: For the specific case of RED, this means that byte-mode queue measurement will often be appropriate although byte-mode drop is strongly deprecated.
NEW: For the specific case of RED, this means that byte-mode queue measurement will often be appropriate but the use of byte-mode drop is very strongly discouraged.

From -10 to -11: Following a further WGLC:

- * Abstract: clarified that advice applies to all AQMs including newer ones
- * Abstract & Intro: changed 'read' to 'detect', because you don't read losses, you detect them.

- * S.1. Introduction: Disambiguated summary of advice on queue measurement.
- * Clarified that the doc deprecates any preference based solely on packet size, it's not only against preferring smaller packets.
- * S.4.1.2. Congestion Measurement without a Queue: Explained that a queue of TXOPs represents a queue into spectrum congested by too many bits.
- * S.5.2: Bit- & Packet-congestible Network: Referred to explanation in S.4.1.2 to make the point that TXOPs are not a primary unit of workload like bits and packets are, even though you get queues of TXOPs.
- * 6. Security: Disambiguated 'bias towards'.
- * 8. Conclusions: Made consistent with recommendation to use time if possible for queue measurement.

From -09 to -10: Following IESG review:

- * Updates 2309: Left header unchanged reflecting eventual IESG consensus [Sean Turner, Pete Resnick].
- * S.1 Intro: This memo adds to the congestion control principles enumerated in BCP 41 [Pete Resnick]
- * Abstract, S.1, S.1.1, s.1.2 Intro, Scoping and Example: Made applicability to all AQMs clearer listing some more example AQMs and explained that we always use RED for examples, but this doesn't mean it's not applicable to other AQMs. [A number of reviewers have described the draft as "about RED"]
- * S.1 & S.2.1 Queue measurement: Explained that the choice between measuring the queue in packets or bytes is only relevant if measuring it in time units is infeasible [So as not to imply that we haven't noticed the advances made by PDPC & CoDel]
- * S.1.1. Terminology: Better explained why hybrid systems congested by both packets and bytes are often designed to be treated as bit-congestible [Richard Barnes].
- * S.2.1. Queue measurement advice: Added examples. Added a counter-example to justify SHOULDs rather than MUSTs. Pointed to S.4.1 for a list of more complicated scenarios. [Benson]

Schliesser, OpsDir]

- * S2.2. Recommendation on Encoding Congestion Notification: Removed SHOULD treat packets equally, leaving only SHOULD NOT drop dependent on packet size, to avoid it sounding like we're saying QoS is not allowed. Pointed to possible app-specific legacy use of byte-mode as a counter-example that prevents us saying MUST NOT. [Pete Resnick]
- * S.2.3. Recommendation on Responding to Congestion: capitalised the two SHOULDs in recommendations for TCP, and gave possible counter-examples. [noticed while dealing with Pete Resnick's point]
- * S2.4. Splitting & Merging: RTCP -> RTP/RTCP [Pete McCann, Gen-ART]
- * S.3.2 Small != Control: many control packets are small -> ...tend to be small [Stephen Farrell]
- * S.3.1 Perverse incentives: Changed transport designers to app developers [Stephen Farrell]
- * S.4.1.1. Fixed Size Packet Buffers: Nearly completely re-written to simplify and to reverse the advice when the underlying resource is bit-congestible, irrespective of whether the buffer consists of fixed-size packet buffers. [Richard Barnes & Benson Schliesser]
- * S.4.2.1.2. Packet Size Bias Regardless of AQM: Largely re-written to reflect the earlier change in advice about fixed-size packet buffers, and to primarily focus on getting rid of tail-drop, not various nuances of tail-drop. [Richard Barnes & Benson Schliesser]
- * Editorial corrections [Tim Bray, AppsDir, Pete McCann, Gen-ART and others]
- * Updated refs (two I-Ds have become RFCs). [Pete McCann]

From -08 to -09: Following WG last call:

- * S.2.1: Made RED-related queue measurement recommendations clearer
- * S.2.3: Added to "Recommendation on Responding to Congestion" to make it clear that we are definitely not saying transports have to equalise bit-rates, just how to do it and not do it, if you

want to.

- * S.3: Clarified motivation sections S.3.3 "Transport-Independent Network" and S.3.5 "Implementation Efficiency"
- * S.3.4: Completely changed motivating argument from "Scaling Congestion Control with Packet Size" to "Partial Deployment of AQM".

From -07 to -08:

- * Altered abstract to say it provides best current practice and highlight that it updates RFC2309
- * Added null IANA section
- * Updated refs

From -06 to -07:

- * A mix-up with the corollaries and their naming in 2.1 to 2.3 fixed.

From -05 to -06:

- * Primarily editorial fixes.

From -04 to -05:

- * Changed from Informational to BCP and highlighted non-normative sections and appendices
- * Removed language about consensus
- * Added "Example Comparing Packet-Mode Drop and Byte-Mode Drop"
- * Arranged "Motivating Arguments" into a more logical order and completely rewrote "Transport-Independent Network" & "Scaling Congestion Control with Packet Size" arguments. Removed "Why Now?"
- * Clarified applicability of certain recommendations
- * Shifted vendor survey to an Appendix
- * Cut down "Outstanding Issues and Next Steps"

- * Re-drafted the start of the conclusions to highlight the three distinct areas of concern
- * Completely re-wrote appendices
- * Editorial corrections throughout.

From -03 to -04:

- * Reordered Sections 2 and 3, and some clarifications here and there based on feedback from Colin Perkins and Mirja Kuehlewind.

From -02 to -03 (this version)

- * Structural changes:
 - + Split off text at end of "Scaling Congestion Control with Packet Size" into new section "Transport-Independent Network"
 - + Shifted "Recommendations" straight after "Motivating Arguments" and added "Conclusions" at end to reinforce Recommendations
 - + Added more internal structure to Recommendations, so that recommendations specific to RED or to TCP are just corollaries of a more general recommendation, rather than being listed as a separate recommendation.
 - + Renamed "State of the Art" as "Critical Survey of Existing Advice" and retitled a number of subsections with more descriptive titles.
 - + Split end of "Congestion Coding: Summary of Status" into a new subsection called "RED Implementation Status".
 - + Removed text that had been in the Appendix "Congestion Notification Definition: Further Justification".
- * Reordered the intro text a little.
- * Made it clearer when advice being reported is deprecated and when it is not.
- * Described AQM as in network equipment, rather than saying "at the network layer" (to side-step controversy over whether functions like AQM are in the transport layer but in network

equipment).

- * Minor improvements to clarity throughout

From -01 to -02:

- * Restructured the whole document for (hopefully) easier reading and clarity. The concrete recommendation, in RFC2119 language, is now in Section 8.

From -00 to -01:

- * Minor clarifications throughout and updated references

From briscoe-byte-pkt-mark-02 to ietf-byte-pkt-congest-00:

- * Added note on relationship to existing RFCs
- * Posed the question of whether packet-congestion could become common and deferred it to the IRTF ICCRG. Added ref to the dual-resource queue (DRQ) proposal.
- * Changed PCN references from the PCN charter & architecture to the PCN marking behaviour draft most likely to imminently become the standards track WG item.

From -01 to -02:

- * Abstract reorganised to align with clearer separation of issue in the memo.
- * Introduction reorganised with motivating arguments removed to new Section 3.
- * Clarified avoiding lock-out of large packets is not the main or only motivation for RED.
- * Mentioned choice of drop or marking explicitly throughout, rather than trying to coin a word to mean either.
- * Generalised the discussion throughout to any packet forwarding function on any network equipment, not just routers.
- * Clarified the last point about why this is a good time to sort out this issue: because it will be hard / impossible to design new transports unless we decide whether the network or the transport is allowing for packet size.

- * Added statement explaining the horizon of the memo is long term, but with short term expediency in mind.
- * Added material on scaling congestion control with packet size (Section 3.4).
- * Separated out issue of normalising TCP's bit rate from issue of preference to control packets (Section 3.2).
- * Divided up Congestion Measurement section for clarity, including new material on fixed size packet buffers and buffer carving (Section 4.1.1 & Section 4.2.1) and on congestion measurement in wireless link technologies without queues (Section 4.1.2).
- * Added section on 'Making Transports Robust against Control Packet Losses' (Section 4.2.3) with existing & new material included.
- * Added tabulated results of vendor survey on byte-mode drop variant of RED (Table 3).

From -00 to -01:

- * Clarified applicability to drop as well as ECN.
- * Highlighted DoS vulnerability.
- * Emphasised that drop-tail suffers from similar problems to byte-mode drop, so only byte-mode drop should be turned off, not RED itself.
- * Clarified the original apparent motivations for recommending byte-mode drop included protecting SYN's and pure ACK's more than equalising the bit rates of TCP's with different segment sizes. Removed some conjectured motivations.
- * Added support for updates to TCP in progress (ackcc & ecn-syn-ack).
- * Updated survey results with newly arrived data.
- * Pulled all recommendations together into the conclusions.
- * Moved some detailed points into two additional appendices and a note.

* Considerable clarifications throughout.

* Updated references

Authors' Addresses

Bob Briscoe
BT
B54/77, Adastral Park
Martlesham Heath
Ipswich IP5 3RE
UK

Phone: +44 1473 645196
EMail: bob.briscoe@bt.com
URI: <http://bobbriscoe.net/>

Jukka Manner
Aalto University
Department of Communications and Networking (Comnet)
P.O. Box 13000
FIN-00076 Aalto
Finland

Phone: +358 9 470 22481
EMail: jukka.manner@aalto.fi
URI: <http://www.netlab.tkk.fi/~jmanner/>

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: 28 April 2022

R. R. Stewart
Netflix, Inc.
M. Tüxen
I. Rüngeler
Münster Univ. of Appl. Sciences
25 October 2021

Stream Control Transmission Protocol (SCTP) Network Address Translation
Support
draft-ietf-tsvwg-natsupp-23

Abstract

The Stream Control Transmission Protocol (SCTP) provides a reliable communications channel between two end-hosts in many ways similar to the Transmission Control Protocol (TCP). With the widespread deployment of Network Address Translators (NAT), specialized code has been added to NAT functions for TCP that allows multiple hosts to reside behind a NAT function and yet share a single IPv4 address, even when two hosts (behind a NAT function) choose the same port numbers for their connection. This additional code is sometimes classified as Network Address and Port Translation (NAPT).

This document describes the protocol extensions needed for the SCTP endpoints and the mechanisms for NAT functions necessary to provide similar features of NAPT in the single point and multipoint traversal scenario.

Finally, a YANG module for SCTP NAT is defined.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

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

This Internet-Draft will expire on 28 April 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 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	3
2. Conventions	5
3. Terminology	5
4. Motivation and Overview	6
4.1. SCTP NAT Traversal Scenarios	6
4.1.1. Single Point Traversal	7
4.1.2. Multipoint Traversal	7
4.2. Limitations of Classical NAPT for SCTP	8
4.3. The SCTP-Specific Variant of NAT	8
5. Data Formats	13
5.1. Modified Chunks	13
5.1.1. Extended ABORT Chunk	13
5.1.2. Extended ERROR Chunk	14
5.2. New Error Causes	14
5.2.1. VTag and Port Number Collision Error Cause	14
5.2.2. Missing State Error Cause	15
5.2.3. Port Number Collision Error Cause	15
5.3. New Parameters	16
5.3.1. Disable Restart Parameter	16
5.3.2. VTags Parameter	17
6. Procedures for SCTP Endpoints and NAT Functions	18
6.1. Association Setup Considerations for Endpoints	19
6.2. Handling of Internal Port Number and Verification Tag Collisions	19
6.2.1. NAT Function Considerations	19
6.2.2. Endpoint Considerations	20
6.3. Handling of Internal Port Number Collisions	20
6.3.1. NAT Function Considerations	20
6.3.2. Endpoint Considerations	21
6.4. Handling of Missing State	21
6.4.1. NAT Function Considerations	22
6.4.2. Endpoint Considerations	22

6.5.	Handling of Fragmented SCTP Packets by NAT Functions . .	24
6.6.	Multi Point Traversal Considerations for Endpoints . . .	24
7.	SCTP NAT YANG Module	24
7.1.	Tree Structure	24
7.2.	YANG Module	25
8.	Various Examples of NAT Traversals	27
8.1.	Single-homed Client to Single-homed Server	28
8.2.	Single-homed Client to Multi-homed Server	30
8.3.	Multihomed Client and Server	32
8.4.	NAT Function Loses Its State	35
8.5.	Peer-to-Peer Communications	37
9.	Socket API Considerations	42
9.1.	Get or Set the NAT Friendliness (SCTP_NAT_FRIENDLY) . . .	43
10.	IANA Considerations	43
10.1.	New Chunk Flags for Two Existing Chunk Types	43
10.2.	Three New Error Causes	45
10.3.	Two New Chunk Parameter Types	46
10.4.	One New URI	46
10.5.	One New YANG Module	46
11.	Security Considerations	46
12.	Normative References	47
13.	Informative References	48
	Acknowledgments	51
	Authors' Addresses	51

1. Introduction

Stream Control Transmission Protocol (SCTP) [RFC4960] provides a reliable communications channel between two end-hosts in many ways similar to TCP [RFC0793]. With the widespread deployment of Network Address Translators (NAT), specialized code has been added to NAT functions for TCP that allows multiple hosts to reside behind a NAT function using private-use addresses (see [RFC6890]) and yet share a single IPv4 address, even when two hosts (behind a NAT function) choose the same port numbers for their connection. This additional code is sometimes classified as Network Address and Port Translation (NAPT). Please note that this document focuses on the case where the NAT function maps a single or multiple internal addresses to a single external address and vice versa.

To date, specialized code for SCTP has not yet been added to most NAT functions so that only a translation of IP addresses is supported. The end result of this is that only one SCTP-capable host can successfully operate behind such a NAT function and this host can only be single-homed. The only alternative for supporting legacy NAT functions is to use UDP encapsulation as specified in [RFC6951].

The NAT function in the document refers to NAPT functions described in Section 2.2 of [RFC3022], NAT64 [RFC6146], or DS-Lite AFTR [RFC6333].

This document specifies procedures allowing a NAT function to support SCTP by providing similar features to those provided by a NAPT for TCP (see [RFC5382] and [RFC7857]), UDP (see [RFC4787] and [RFC7857]), and ICMP (see [RFC5508] and [RFC7857]). This document also specifies a set of data formats for SCTP packets and a set of SCTP endpoint procedures to support NAT traversal. An SCTP implementation supporting these procedures can assure that in both single-homed and multi-homed cases a NAT function will maintain the appropriate state without the NAT function needing to change port numbers.

It is possible and desirable to make these changes for a number of reasons:

- * It is desirable for SCTP internal end-hosts on multiple platforms to be able to share a NAT function's external IP address in the same way that a TCP session can use a NAT function.
- * If a NAT function does not need to change any data within an SCTP packet, it will reduce the processing burden of NAT'ing SCTP by not needing to execute the CRC32c checksum used by SCTP.
- * Not having to touch the IP payload makes the processing of ICMP messages by NAT functions easier.

An SCTP-aware NAT function will need to follow these procedures for generating appropriate SCTP packet formats.

When considering SCTP-aware NAT it is possible to have multiple levels of support. At each level, the Internal Host, Remote Host, and NAT function does or does not support the procedures described in this document. The following table illustrates the results of the various combinations of support and if communications can occur between two endpoints.

Internal Host	NAT Function	Remote Host	Communication
Support	Support	Support	Yes
Support	Support	No Support	Limited
Support	No Support	Support	None
Support	No Support	No Support	None
No Support	Support	Support	Limited
No Support	Support	No Support	Limited
No Support	No Support	Support	None
No Support	No Support	No Support	None

Table 1: Communication possibilities

From the table it can be seen that no communication can occur when a NAT function does not support SCTP-aware NAT. This assumes that the NAT function does not handle SCTP packets at all and all SCTP packets sent from behind a NAT function are discarded by the NAT function. In some cases, where the NAT function supports SCTP-aware NAT, but one of the two hosts does not support the feature, communication can possibly occur in a limited way. For example, only one host can have a connection when a collision case occurs.

2. Conventions

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.

3. Terminology

This document uses the following terms, which are depicted in Figure 1. Familiarity with the terminology used in [RFC4960] and [RFC5061] is assumed.

Internal-Address (Int-Addr)

An internal address that is known to the internal host.

Internal-Port (Int-Port)

The port number that is in use by the host holding the Internal-Address.

Internal-VTag (Int-VTag)

The SCTP Verification Tag (VTag) (see Section 3.1 of [RFC4960]) that the internal host has chosen for an association. The VTag is a unique 32-bit tag that accompanies any incoming SCTP packet for this association to the Internal-Address.

Remote-Address (Rem-Addr)

The address that an internal host is attempting to contact.

Remote-Port (Rem-Port)

The port number used by the host holding the Remote-Address.

Remote-VTag (Rem-VTag)

The Verification Tag (VTag) (see Section 3.1 of [RFC4960]) that the host holding the Remote-Address has chosen for an association. The VTag is a unique 32-bit tag that accompanies any outgoing SCTP packet for this association to the Remote-Address.

External-Address (Ext-Addr)

An external address assigned to the NAT function, that it uses as a source address when sending packets towards a Remote-Address.

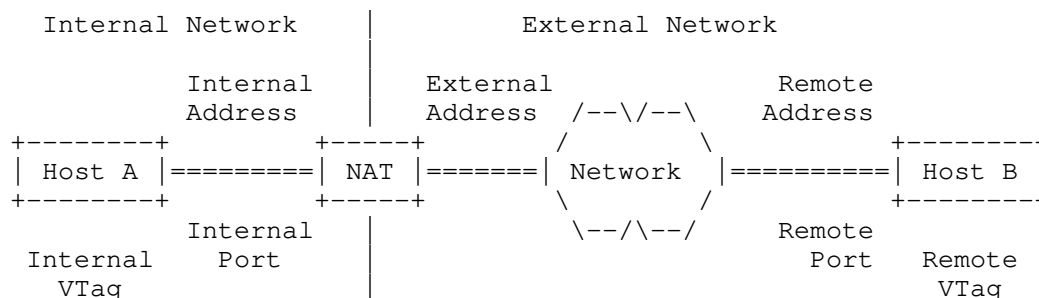


Figure 1: Basic Network Setup

4. Motivation and Overview

4.1. SCTP NAT Traversal Scenarios

This section defines the notion of single and multipoint NAT traversal.

4.1.1. Single Point Traversal

In this case, all packets in the SCTP association go through a single NAT function, as shown in Figure 2.

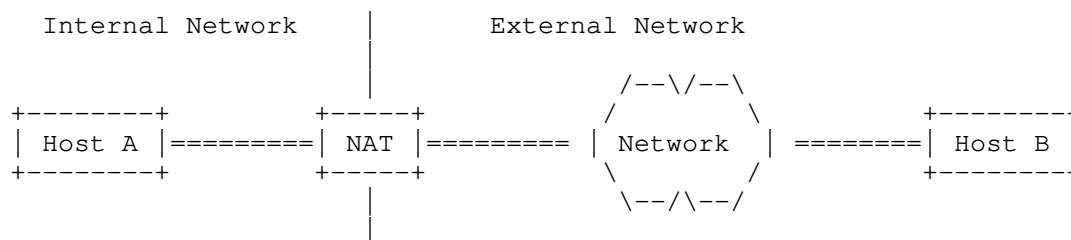


Figure 2: Single NAT Function Scenario

A variation of this case is shown in Figure 3, i.e., multiple NAT functions in the forwarding path between two endpoints.

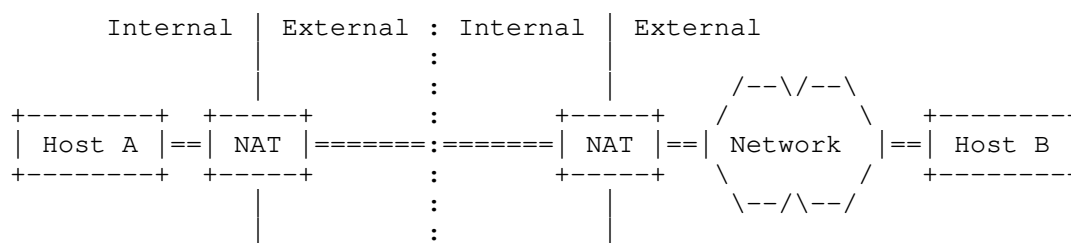


Figure 3: Serial NAT Functions Scenario

Although one of the main benefits of SCTP multi-homing is redundant paths, in the single point traversal scenario the NAT function represents a single point of failure in the path of the SCTP multi-homed association. However, the rest of the path can still benefit from path diversity provided by SCTP multi-homing.

The two SCTP endpoints in this case can be either single-homed or multi-homed. However, the important thing is that the NAT function in this case sees all the packets of the SCTP association.

4.1.2. Multipoint Traversal

This case involves multiple NAT functions and each NAT function only sees some of the packets in the SCTP association. An example is shown in Figure 4.

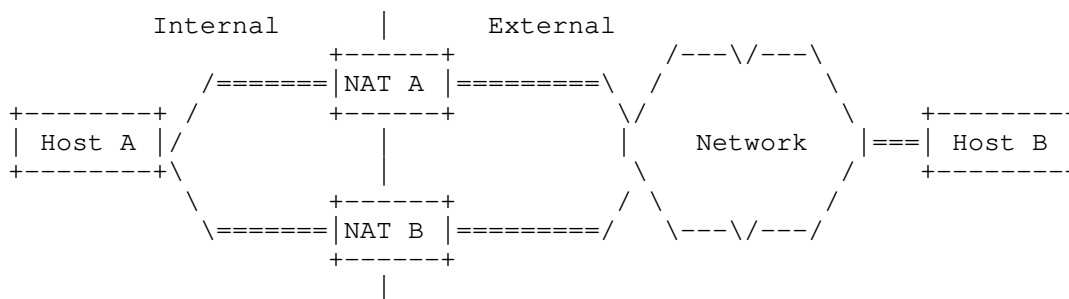


Figure 4: Parallel NAT Functions Scenario

This case does not apply to a single-homed SCTP association (i.e., both endpoints in the association use only one IP address). The advantage here is that the existence of multiple NAT traversal points can preserve the path diversity of a multi-homed association for the entire path. This in turn can improve the robustness of the communication.

4.2. Limitations of Classical NAPT for SCTP

Using classical NAPT possibly results in changing one of the SCTP port numbers during the processing, which requires the recomputation of the transport layer checksum by the NAPT function. Whereas for UDP and TCP this can be done very efficiently, for SCTP the checksum (CRC32c) over the entire packet needs to be recomputed (see Appendix B of [RFC4960] for details of the CRC32c computation). This would considerably add to the NAT computational burden, however hardware support can mitigate this in some implementations.

An SCTP endpoint can have multiple addresses but only has a single port number to use. To make multipoint traversal work, all the NAT functions involved need to recognize the packets they see as belonging to the same SCTP association and perform port number translation in a consistent way. One possible way of doing this is to use a pre-defined table of port numbers and addresses configured within each NAT function. Other mechanisms could make use of NAT to NAT communication. Such mechanisms have not been deployed on a wide scale base and thus are not a preferred solution. Therefore an SCTP variant of NAT function has been developed (see Section 4.3).

4.3. The SCTP-Specific Variant of NAT

In this section it is allowed that there are multiple SCTP capable hosts behind a NAT function that share one External-Address. Furthermore, this section focuses on the single point traversal scenario (see Section 4.1.1).

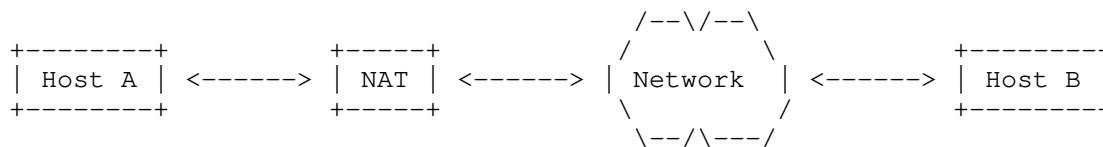
The modification of outgoing SCTP packets sent from an internal host is simple: the source address of the packets has to be replaced with the External-Address. It might also be necessary to establish some state in the NAT function to later handle incoming packets.

Typically, the NAT function has to maintain a NAT binding table of Internal-VTag, Internal-Port, Remote-VTag, Remote-Port, Internal-Address, and whether the restart procedure is disabled or not. An entry in that NAT binding table is called a NAT-State control block. The function Create() obtains the just mentioned parameters and returns a NAT-State control block. A NAT function MAY allow creating NAT-State control blocks via a management interface.

For SCTP packets coming from the external realm of the NAT function the destination address of the packets has to be replaced with the Internal-Address of the host to which the packet has to be delivered, if a NAT state entry is found. The lookup of the Internal-Address is based on the Remote-VTag, Remote-Port, Internal-VTag and the Internal-Port.

The entries in the NAT binding table need to fulfill some uniqueness conditions. There can not be more than one entry NAT binding table with the same pair of Internal-Port and Remote-Port. This rule can be relaxed, if all NAT binding table entries with the same Internal-Port and Remote-Port have the support for the restart procedure disabled (see Section 5.3.1). In this case there can not be no more than one entry with the same Internal-Port, Remote-Port and Remote-VTag and no more than one NAT binding table entry with the same Internal-Port, Remote-Port, and Int-VTag.

The processing of outgoing SCTP packets containing an INIT chunk is illustrated in the following figure. This scenario is valid for all message flows in this section.



```

INIT[Initiate-Tag]
Int-Addr:Int-Port -----> Rem-Addr:Rem-Port
Rem-VTag=0

Create(Initiate-Tag, Int-Port, 0, Rem-Port, Int-Addr,
      IsRestartDisabled)
Returns(NAT-State control block)

```

Translate To:

```

INIT[Initiate-Tag]
Ext-Addr:Int-Port -----> Rem-Addr:Rem-Port
Rem-VTag=0

```

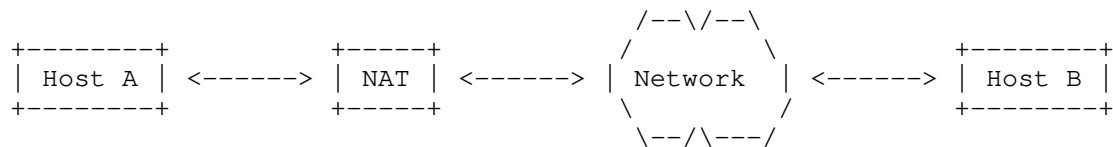
Normally a NAT binding table entry will be created.

However, it is possible that there is already a NAT binding table entry with the same Remote-Port, Internal-Port, and Internal-VTag but different Internal-Address and the restart procedure is disabled. In this case the packet containing the INIT chunk MUST be dropped by the NAT and a packet containing an ABORT chunk SHOULD be sent to the SCTP host that originated the packet with the M bit set and 'VTag and Port Number Collision' error cause (see Section 5.1.1 for the format). The source address of the packet containing the ABORT chunk MUST be the destination address of the packet containing the INIT chunk.

If an outgoing SCTP packet contains an INIT or ASCONF chunk and a matching NAT binding table entry is found, the packet is processed as a normal outgoing packet.

It is also possible that a NAT binding table entry with the same Remote-Port and Internal-Port exists without an Internal-VTag conflict but there exists a NAT binding table entry with the same port numbers but a different Internal-Address and the restart procedure is not disabled. In such a case the packet containing the INIT chunk MUST be dropped by the NAT function and a packet containing an ABORT chunk SHOULD be sent to the SCTP host that originated the packet with the M bit set and 'Port Number Collision' error cause (see Section 5.1.1 for the format).

The processing of outgoing SCTP packets containing no INIT chunks is described in the following figure.

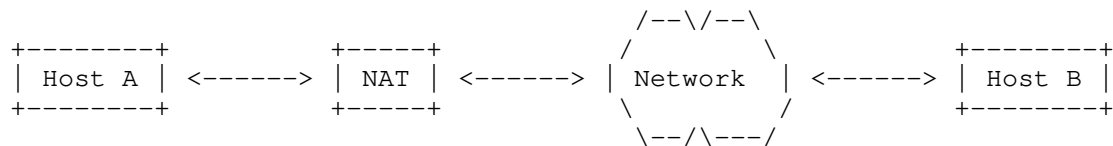


Int-Addr:Int-Port -----> Rem-Addr:Rem-Port
 Rem-VTag

Translate To:

Ext-Addr:Int-Port -----> Rem-Addr:Rem-Port
 Rem-VTag

The processing of incoming SCTP packets containing an INIT ACK chunk is illustrated in the following figure. The Lookup() function has as input the Internal-VTag, Internal-Port, Remote-VTag, and Remote-Port. It returns the corresponding entry of the NAT binding table and updates the Remote-VTag by substituting it with the value of the Initiate-Tag of the INIT ACK chunk. The wildcard character signifies that the parameter's value is not considered in the Lookup() function or changed in the Update() function, respectively.



INIT ACK[Initiate-Tag]
 Ext-Addr:Int-Port <---- Rem-Addr:Rem-Port
 Int-VTag

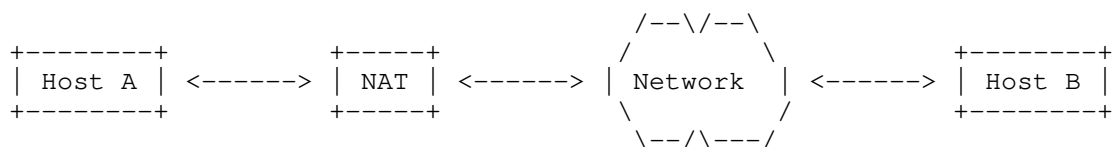
Lookup(Int-VTag, Int-Port, *, Rem-Port)
 Update(*, *, Initiate-Tag, *)

Returns(NAT-State control block containing Int-Addr)

INIT ACK[Initiate-Tag]
 Int-Addr:Int-Port <----- Rem-Addr:Rem-Port
 Int-VTag

In the case where the Lookup function fails because it does not find an entry, the SCTP packet is dropped. If it succeeds, the Update routine inserts the Remote-VTag (the Initiate-Tag of the INIT ACK chunk) in the NAT-State control block.

The processing of incoming SCTP packets containing an ABORT or SHUTDOWN COMPLETE chunk with the T bit set is illustrated in the following figure.



Ext-Addr:Int-Port <----- Rem-Addr:Rem-Port
Rem-VTag

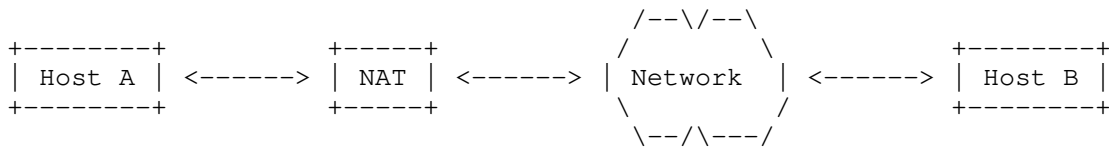
Lookup(*, Int-Port, Rem-VTag, Rem-Port)

Returns (NAT-State control block containing Int-Addr)

Int-Addr:Int-Port <----- Rem-Addr:Rem-Port
Rem-VTag

For an incoming packet containing an INIT chunk a table lookup is made only based on the addresses and port numbers. If an entry with a Remote-VTag of zero is found, it is considered a match and the Remote-VTag is updated. If an entry with a non-matching Remote-VTag is found or no entry is found, the incoming packet is silently dropped. If an entry with a matching Remote-VTag is found, the incoming packet is forwarded. This allows the handling of INIT collision through NAT functions.

The processing of other incoming SCTP packets is described in the following figure.



Ext-Addr: Int-Port <----- Rem-Addr: Rem-Port
Int-VTag

Lookup(Int-VTag, Int-Port, *, Rem-Port)

Returns(NAT-State control block containing Internal-Address)

Int-Addr: Int-Port <----- Rem-Addr: Rem-Port
Int-VTag

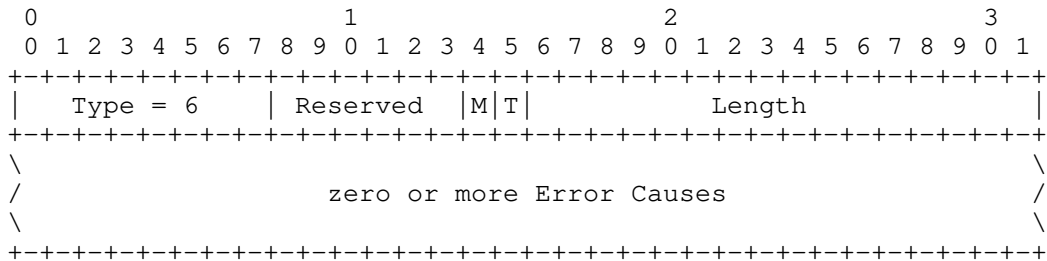
5. Data Formats

This section defines the formats used to support NAT traversal. Section 5.1 and Section 5.2 describe chunks and error causes sent by NAT functions and received by SCTP endpoints. Section 5.3 describes parameters sent by SCTP endpoints and used by NAT functions and SCTP endpoints.

5.1. Modified Chunks

This section presents existing chunks defined in [RFC4960] for which additional flags are specified by this document.

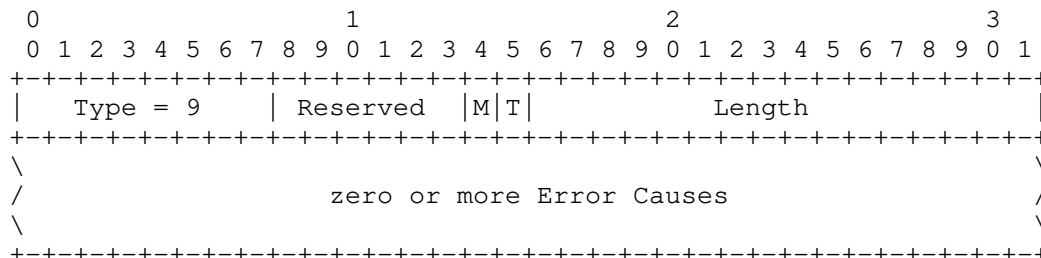
5.1.1. Extended ABORT Chunk



The ABORT chunk is extended to add the new 'M bit'. The M bit indicates to the receiver of the ABORT chunk that the chunk was not generated by the peer SCTP endpoint, but instead by a middle box (e.g., NAT).

[NOTE to RFC-Editor: Assignment of M bit to be confirmed by IANA.]

5.1.2. Extended ERROR Chunk



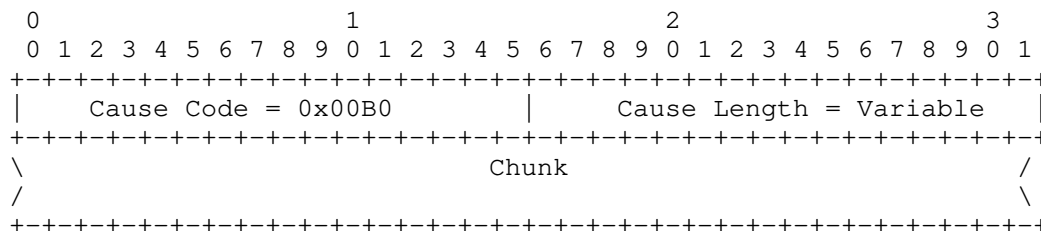
The ERROR chunk defined in [RFC4960] is extended to add the new 'M bit'. The M bit indicates to the receiver of the ERROR chunk that the chunk was not generated by the peer SCTP endpoint, but instead by a middle box.

[NOTE to RFC-Editor: Assignment of M bit to be confirmed by IANA.]

5.2. New Error Causes

This section defines the new error causes added by this document.

5.2.1. VTag and Port Number Collision Error Cause



Cause Code: 2 bytes (unsigned integer)

This field holds the IANA defined cause code for the 'VTag and Port Number Collision' Error Cause. IANA is requested to assign the value 0x00B0 for this cause code.

Cause Length: 2 bytes (unsigned integer)

This field holds the length in bytes of the error cause. The value MUST be the length of the Cause-Specific Information plus 4.

Chunk: variable length

The Cause-Specific Information is filled with the chunk that caused this error. This can be an INIT, INIT ACK, or ASCONF chunk. Note that if the entire chunk will not fit in the ERROR chunk or ABORT chunk being sent then the bytes that do not fit are truncated.

[NOTE to RFC-Editor: Assignment of cause code to be confirmed by IANA.]

5.2.2. Missing State Error Cause

0										1										2										3									
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1								
Cause Code = 0x00B1										Cause Length = Variable																													
Original Packet																																							

Cause Code: 2 bytes (unsigned integer)

This field holds the IANA defined cause code for the 'Missing State' Error Cause. IANA is requested to assign the value 0x00B1 for this cause code.

Cause Length: 2 bytes (unsigned integer)

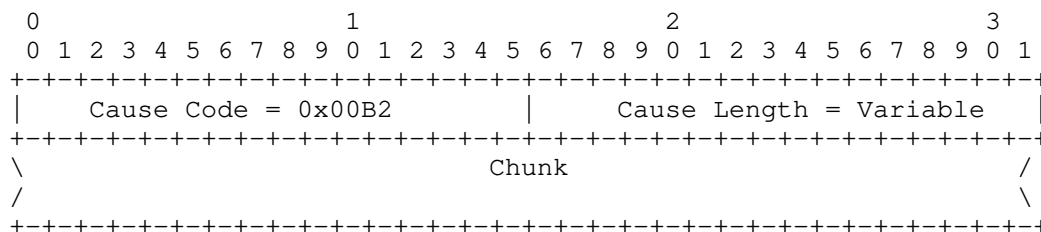
This field holds the length in bytes of the error cause. The value MUST be the length of the Cause-Specific Information plus 4.

Original Packet: variable length

The Cause-Specific Information is filled with the IPv4 or IPv6 packet that caused this error. The IPv4 or IPv6 header MUST be included. Note that if the packet will not fit in the ERROR chunk or ABORT chunk being sent then the bytes that do not fit are truncated.

[NOTE to RFC-Editor: Assignment of cause code to be confirmed by IANA.]

5.2.3. Port Number Collision Error Cause



Cause Code: 2 bytes (unsigned integer)

This field holds the IANA defined cause code for the 'Port Number Collision' Error Cause. IANA is requested to assign the value 0x00B2 for this cause code.

Cause Length: 2 bytes (unsigned integer)

This field holds the length in bytes of the error cause. The value MUST be the length of the Cause-Specific Information plus 4.

Chunk: variable length

The Cause-Specific Information is filled with the chunk that caused this error. This can be an INIT, INIT ACK, or ASCONF chunk. Note that if the entire chunk will not fit in the ERROR chunk or ABORT chunk being sent then the bytes that do not fit are truncated.

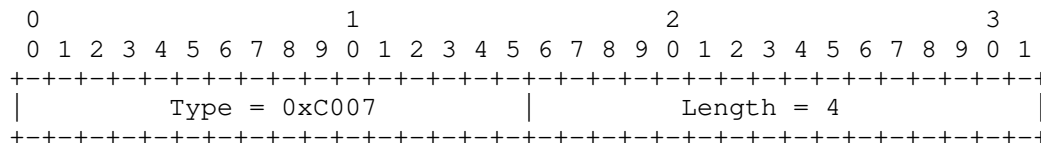
[NOTE to RFC-Editor: Assignment of cause code to be confirmed by IANA.]

5.3. New Parameters

This section defines new parameters and their valid appearance defined by this document.

5.3.1. Disable Restart Parameter

This parameter is used to indicate that the restart procedure is requested to be disabled. Both endpoints of an association MUST include this parameter in the INIT chunk and INIT ACK chunk when establishing an association and MUST include it in the ASCONF chunk when adding an address to successfully disable the restart procedure.



Parameter Type: 2 bytes (unsigned integer)

This field holds the IANA defined parameter type for the Disable Restart Parameter. IANA is requested to assign the value 0xC007 for this parameter type.

Parameter Length: 2 bytes (unsigned integer)

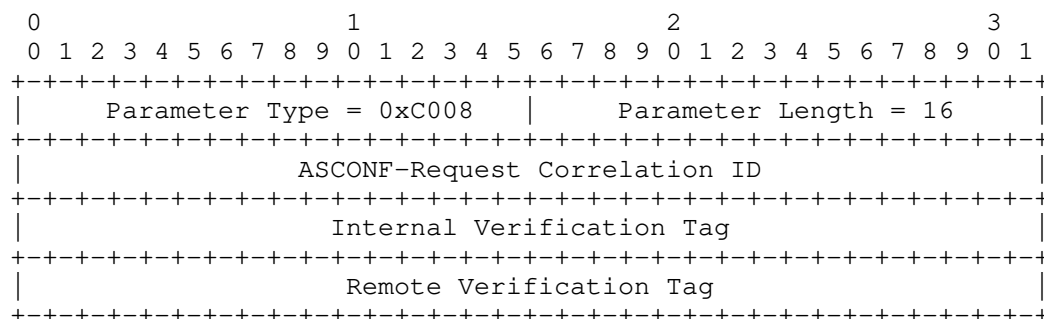
This field holds the length in bytes of the parameter. The value MUST be 4.

[NOTE to RFC-Editor: Assignment of parameter type to be confirmed by IANA.]

The Disable Restart Parameter MAY appear in INIT, INIT ACK and ASCONF chunks and MUST NOT appear in any other chunk.

5.3.2. VTags Parameter

This parameter is used to help a NAT function to recover from state loss.



Parameter Type: 2 bytes (unsigned integer)

This field holds the IANA defined parameter type for the VTags Parameter. IANA is requested to assign the value 0xC008 for this parameter type.

Parameter Length: 2 bytes (unsigned integer)

This field holds the length in bytes of the parameter. The value MUST be 16.

ASCONF-Request Correlation ID: 4 bytes (unsigned integer)

This is an opaque integer assigned by the sender to identify each request parameter. The receiver of the ASCONF Chunk will copy this 32-bit value into the ASCONF Response Correlation ID field of the ASCONF ACK response parameter. The sender of the packet containing the ASCONF chunk can use this same value in the ASCONF ACK chunk to find which request the response is for. The receiver MUST NOT change the value of the ASCONF-Request Correlation ID.

Internal Verification Tag: 4 bytes (unsigned integer)

The Verification Tag that the internal host has chosen for the association. The Verification Tag is a unique 32-bit tag that accompanies any incoming SCTP packet for this association to the Internal-Address.

Remote Verification Tag: 4 bytes (unsigned integer)

The Verification Tag that the host holding the Remote-Address has chosen for the association. The VTag is a unique 32-bit tag that accompanies any outgoing SCTP packet for this association to the Remote-Address.

[NOTE to RFC-Editor: Assignment of parameter type to be confirmed by IANA.]

The VTags Parameter MAY appear in ASCONF chunks and MUST NOT appear in any other chunk.

6. Procedures for SCTP Endpoints and NAT Functions

If an SCTP endpoint is behind an SCTP-aware NAT, a number of problems can arise as it tries to communicate with its peers:

- * IP addresses can not be included in the SCTP packet. This is discussed in Section 6.1.
- * More than one host behind a NAT function could select the same VTag and source port number when communicating with the same peer server. This creates a situation where the NAT function will not be able to tell the two associations apart. This situation is discussed in Section 6.2.
- * If an SCTP endpoint is a server communicating with multiple peers and the peers are behind the same NAT function, then these peers cannot be distinguished by the server. This case is discussed in Section 6.3.
- * A restart of a NAT function during a conversation could cause a loss of its state. This problem and its solution is discussed in Section 6.4.
- * NAT functions need to deal with SCTP packets being fragmented at the IP layer. This is discussed in Section 6.5.
- * An SCTP endpoint can be behind two NAT functions in parallel providing redundancy. The method to set up this scenario is discussed in Section 6.6.

The mechanisms to solve these problems require additional chunks and parameters, defined in this document, and modified handling procedures from those specified in [RFC4960] as described below.

6.1. Association Setup Considerations for Endpoints

The association setup procedure defined in [RFC4960] allows multi-homed SCTP endpoints to exchange its IP-addresses by using IPv4 or IPv6 address parameters in the INIT and INIT ACK chunks. However, this does not work when NAT functions are present.

Every association setup from a host behind a NAT function MUST NOT use multiple internal addresses. The INIT chunk MUST NOT contain an IPv4 Address parameter, IPv6 Address parameter, or Supported Address Types parameter. The INIT ACK chunk MUST NOT contain any IPv4 Address parameter or IPv6 Address parameter using non-global addresses. The INIT chunk and the INIT ACK chunk MUST NOT contain any Host Name parameters.

If the association is intended to be finally multi-homed, the procedure in Section 6.6 MUST be used.

The INIT and INIT ACK chunk SHOULD contain the Disable Restart parameter defined in Section 5.3.1.

6.2. Handling of Internal Port Number and Verification Tag Collisions

Consider the case where two hosts in the Internal-Address space want to set up an SCTP association with the same service provided by some remote hosts. This means that the Remote-Port is the same. If they both choose the same Internal-Port and Internal-VTag, the NAT function cannot distinguish between incoming packets anymore. However, this is unlikely. The Internal-VTags are chosen at random and if the Internal-Ports are also chosen from the ephemeral port range at random (see [RFC6056]) this gives a 46-bit random number that has to match.

The same can happen with the Remote-VTag when a packet containing an INIT ACK chunk or an ASCONF chunk is processed by the NAT function.

6.2.1. NAT Function Considerations

If the NAT function detects a collision of internal port numbers and verification tags, it SHOULD send a packet containing an ABORT chunk with the M bit set if the collision is triggered by a packet containing an INIT or INIT ACK chunk. If such a collision is triggered by a packet containing an ASCONF chunk, it SHOULD send a packet containing an ERROR chunk with the M bit. The M bit is a new

bit defined by this document to express to SCTP that the source of this packet is a "middle" box, not the peer SCTP endpoint (see Section 5.1.1). If a packet containing an INIT ACK chunk triggers the collision, the corresponding packet containing the ABORT chunk MUST contain the same source and destination address and port numbers as the packet containing the INIT ACK chunk. If a packet containing an INIT chunk or an ASCONF chunk, the source and destination address and port numbers MUST be swapped.

The sender of the packet containing an ERROR or ABORT chunk MUST include the error cause with cause code 'VTag and Port Number Collision' (see Section 5.2.1).

6.2.2. Endpoint Considerations

The sender of the packet containing the INIT chunk or the receiver of a packet containing the INIT ACK chunk, upon reception of a packet containing an ABORT chunk with M bit set and the appropriate error cause code for colliding NAT binding table state is included, SHOULD reinitiate the association setup procedure after choosing a new initiate tag, if the association is in COOKIE-WAIT state. In any other state, the SCTP endpoint MUST NOT respond.

The sender of the packet containing the ASCONF chunk, upon reception of a packet containing an ERROR chunk with M bit set, MUST stop adding the path to the association.

6.3. Handling of Internal Port Number Collisions

When two SCTP hosts are behind an SCTP-aware NAT it is possible that two SCTP hosts in the Internal-Address space will want to set up an SCTP association with the same server running on the same remote host. If the two hosts choose the same internal port, this is considered an internal port number collision.

For the NAT function, appropriate tracking can be performed by assuring that the VTags are unique between the two hosts.

6.3.1. NAT Function Considerations

The NAT function, when processing the packet containing the INIT ACK chunk, SHOULD note in its NAT binding table if the association supports the disable restart extension. This note is used when establishing future associations (i.e. when processing a packet containing an INIT chunk from an internal host) to decide if the connection can be allowed. The NAT function does the following when processing a packet containing an INIT chunk:

- * If the packet containing the INIT chunk is originating from an internal port to a remote port for which the NAT function has no matching NAT binding table entry, it MUST allow the packet containing the INIT chunk creating an NAT binding table entry.
- * If the packet containing the INIT chunk matches an existing NAT binding table entry, it MUST validate that the disable restart feature is supported and, if it does, allow the packet containing the INIT chunk to be forwarded.
- * If the disable restart feature is not supported, the NAT function SHOULD send a packet containing an ABORT chunk with the M bit set.

The 'Port Number Collision' error cause (see Section 5.2.3) MUST be included in the ABORT chunk sent in response to the packet containing an INIT chunk.

If the collision is triggered by a packet containing an ASCONF chunk, a packet containing an ERROR chunk with the 'Port Number Collision' error cause SHOULD be sent in response to the packet containing the ASCONF chunk.

6.3.2. Endpoint Considerations

For the remote SCTP server this means that the Remote-Port and the Remote-Address are the same. If they both have chosen the same Internal-Port the server cannot distinguish between both associations based on the address and port numbers. For the server it looks like the association is being restarted. To overcome this limitation the client sends a Disable Restart parameter in the INIT chunk.

When the server receives this parameter it does the following:

- * It MUST include a Disable Restart parameter in the INIT ACK to inform the client that it will support the feature.
- * It MUST disable the restart procedures defined in [RFC4960] for this association.

Servers that support this feature will need to be capable of maintaining multiple connections to what appears to be the same peer (behind the NAT function) differentiated only by the VTags.

6.4. Handling of Missing State

6.4.1. NAT Function Considerations

If the NAT function receives a packet from the internal network for which the lookup procedure does not find an entry in the NAT binding table, a packet containing an ERROR chunk SHOULD be sent back with the M bit set. The source address of the packet containing the ERROR chunk MUST be the destination address of the packet received from the internal network. The verification tag is reflected and the T bit is set. Such a packet containing an ERROR chunk SHOULD NOT be sent if the received packet contains an ASCONF chunk with the VTags parameter or an ABORT, SHUTDOWN COMPLETE or INIT ACK chunk. A packet containing an ERROR chunk MUST NOT be sent if the received packet contains an ERROR chunk with the M bit set. In any case, the packet SHOULD NOT be forwarded to the remote address.

If the NAT function receives a packet from the internal network for which it has no NAT binding table entry and the packet contains an ASCONF chunk with the VTags parameter, the NAT function MUST update its NAT binding table according to the verification tags in the VTags parameter and, if present, the Disable Restart parameter.

When sending a packet containing an ERROR chunk, the error cause 'Missing State' (see Section 5.2.2) MUST be included and the M bit of the ERROR chunk MUST be set (see Section 5.1.2).

6.4.2. Endpoint Considerations

Upon reception of this packet containing the ERROR chunk by an SCTP endpoint the receiver takes the following actions:

- * It SHOULD validate that the verification tag is reflected by looking at the VTag that would have been included in an outgoing packet. If the validation fails, discard the received packet containing the ERROR chunk.
- * It SHOULD validate that the peer of the SCTP association supports the dynamic address extension. If the validation fails, discard the received packet containing the ERROR chunk.
- * It SHOULD generate a packet containing a new ASCONF chunk containing the VTags parameter (see Section 5.3.2) and the Disable Restart parameter (see Section 5.3.1) if the association is using the disable restart feature. By processing this packet the NAT function can recover the appropriate state. The procedures for generating an ASCONF chunk can be found in [RFC5061].

The peer SCTP endpoint receiving such a packet containing an ASCONF chunk SHOULD add the address and respond with an acknowledgment if the address is new to the association (following all procedures defined in [RFC5061]). If the address is already part of the association, the SCTP endpoint MUST NOT respond with an error, but instead SHOULD respond with a packet containing an ASCONF ACK chunk acknowledging the address and take no action (since the address is already in the association).

Note that it is possible that upon receiving a packet containing an ASCONF chunk containing the VTags parameter the NAT function will realize that it has an 'Internal Port Number and Verification Tag collision'. In such a case the NAT function SHOULD send a packet containing an ERROR chunk with the error cause code set to 'VTag and Port Number Collision' (see Section 5.2.1).

If an SCTP endpoint receives a packet containing an ERROR chunk with 'Internal Port Number and Verification Tag collision' as the error cause and the packet in the Error Chunk contains an ASCONF with the VTags parameter, careful examination of the association is necessary. The endpoint does the following:

- * It MUST validate that the verification tag is reflected by looking at the VTag that would have been included in the outgoing packet. If the validation fails, it MUST discard the packet.
- * It MUST validate that the peer of the SCTP association supports the dynamic address extension. If the peer does not support this extension, it MUST discard the received packet containing the ERROR chunk.
- * If the association is attempting to add an address (i.e. following the procedures in Section 6.6) then the endpoint MUST NOT consider the address part of the association and SHOULD make no further attempt to add the address (i.e. cancel any ASCONF timers and remove any record of the path), since the NAT function has a VTag collision and the association cannot easily create a new VTag (as it would if the error occurred when sending a packet containing an INIT chunk).
- * If the endpoint has no other path, i.e. the procedure was executed due to missing a state in the NAT function, then the endpoint MUST abort the association. This would occur only if the local NAT function restarted and accepted a new association before attempting to repair the missing state (Note that this is no different than what happens to all TCP connections when a NAT function loses its state).

6.5. Handling of Fragmented SCTP Packets by NAT Functions

SCTP minimizes the use of IP-level fragmentation. However, it can happen that using IP-level fragmentation is needed to continue an SCTP association. For example, if the path MTU is reduced and there are still some DATA chunk in flight, which require packets larger than the new path MTU. If IP-level fragmentation can not be used, the SCTP association will be terminated in a non-graceful way. See [RFC8900] for more information about IP fragmentation.

Therefore, a NAT function MUST be able to handle IP-level fragmented SCTP packets. The fragments MAY arrive in any order.

When an SCTP packet can not be forwarded by the NAT function due to MTU issues and the IP header forbids fragmentation, the NAT MUST send back a "Fragmentation needed and DF set" ICMPv4 or PTB ICMPv6 message to the internal host. This allows for a faster recovery from this packet drop.

6.6. Multi Point Traversal Considerations for Endpoints

If a multi-homed SCTP endpoint behind a NAT function connects to a peer, it MUST first set up the association single-homed with only one address causing the first NAT function to populate its state. Then it SHOULD add each IP address using packets containing ASCONF chunks sent via their respective NAT functions. The address used in the Add IP address parameter is the wildcard address (0.0.0.0 or ::0) and the address parameter in the ASCONF chunk SHOULD also contain the VTags parameter and optionally the Disable Restart parameter.

7. SCTP NAT YANG Module

This section defines a YANG module for SCTP NAT.

The terminology for describing YANG data models is defined in [RFC7950]. The meaning of the symbols in tree diagrams is defined in [RFC8340].

7.1. Tree Structure

This module augments NAT YANG module [RFC8512] with SCTP specifics. The module supports both classical SCTP NAT (that is, rewrite port numbers) and SCTP-specific variant where the ports numbers are not altered. The YANG "feature" is used to indicate whether SCTP-specific variant is supported.

The tree structure of the SCTP NAT YANG module is provided below:

```

module: ietf-nat-sctp
  augment /nat:nat/nat:instances/nat:instance
    /nat:policy/nat:timers:
      +--rw sctp-timeout?  uint32
  augment /nat:nat/nat:instances/nat:instance
    /nat:mapping-table/nat:mapping-entry:
      +--rw int-VTag?      uint32 {sctp-nat}?
      +--rw rem-VTag?      uint32 {sctp-nat}?

```

Concretely, the SCTP NAT YANG module augments the NAT YANG module (policy, in particular) with the following:

- * The sctp-timeout is used to control the SCTP inactivity timeout. That is, the time an SCTP mapping will stay active without SCTP packets traversing the NAT. This timeout can be set only for SCTP. Hence, `"/nat:nat/nat:instances/nat:instance/nat:policy/nat:transport-protocols/nat:protocol-id"` MUST be set to `'132'` (SCTP).

In addition, the SCTP NAT YANG module augments the mapping entry with the following parameters defined in Section 3. These parameters apply only for SCTP NAT mapping entries (i.e., `"/nat/instances/instance/mapping-table/mapping-entry/transport-protocol"` MUST be set to `'132'`);

- * The Internal Verification Tag (Int-VTag)
- * The Remote Verification Tag (Rem-VTag)

7.2. YANG Module

```

<CODE BEGINS> file "ietf-nat-sctp@2020-11-02.yang"
module ietf-nat-sctp {
  yang-version 1.1;
  namespace "urn:ietf:params:xml:ns:yang:ietf-nat-sctp";
  prefix nat-sctp;

  import ietf-nat {
    prefix nat;
    reference
      "RFC 8512: A YANG Module for Network Address Translation
       (NAT) and Network Prefix Translation (NPT)";
  }

  organization
    "IETF TSVWG Working Group";
  contact
    "WG Web:  <https://datatracker.ietf.org/wg/tsvwg/>

```


WG List: <mailto:tsvwg@ietf.org>

Author: Mohamed Boucadair
<mailto:mohamed.boucadair@orange.com>;

description

"This module augments NAT YANG module with Stream Control Transmission Protocol (SCTP) specifics. The extension supports both a classical SCTP NAT (that is, rewrite port numbers) and a, SCTP-specific variant where the ports numbers are not altered.

Copyright (c) 2020 IETF Trust and the persons identified as authors of the code. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, is permitted pursuant to, and subject to the license terms contained in, the Simplified BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>).

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.";

revision 2019-11-18 {

description

"Initial revision.";

reference

"RFC XXXX: Stream Control Transmission Protocol (SCTP) Network Address Translation Support";

}

feature sctp-nat {

description

"This feature means that SCTP-specific variant of NAT is supported. That is, avoid rewriting port numbers.";

reference

"Section 4.3 of RFC XXXX.";

}

augment "/nat:nat/nat:instances/nat:instance"

+ "/nat:policy/nat:timers" {

when "/nat:nat/nat:instances/nat:instance"

+ "/nat:policy/nat:transport-protocols"

+ "/nat:protocol-id = 132";

description

"Extends NAT policy with a timeout for SCTP mapping entries.";

```
leaf sctp-timeout {
  type uint32;
  units "seconds";
  description
    "SCTP inactivity timeout. That is, the time an SCTP
    mapping entry will stay active without packets
    traversing the NAT.";
}

augment "/nat:nat/nat:instances/nat:instance"
  + "/nat:mapping-table/nat:mapping-entry" {
  when "nat:transport-protocol = 132";
  if-feature "sctp-nat";
  description
    "Extends the mapping entry with SCTP specifics.";

  leaf int-VTag {
    type uint32;
    description
      "The Internal Verification Tag that the internal
      host has chosen for this communication.";
  }
  leaf rem-VTag {
    type uint32;
    description
      "The Remote Verification Tag that the remote
      peer has chosen for this communication.";
  }
}
}
<CODE ENDS>
```

8. Various Examples of NAT Traversals

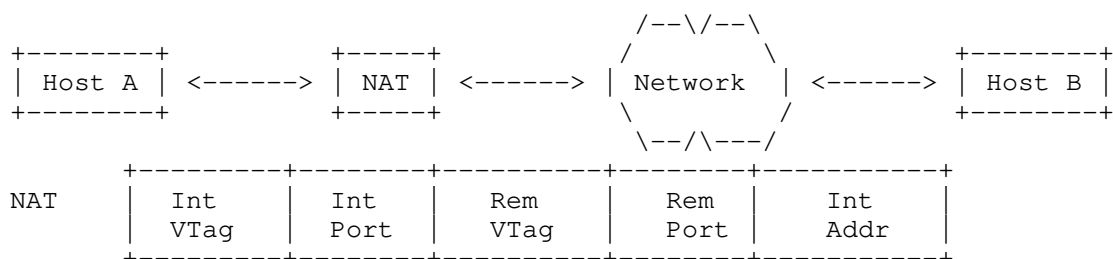
Please note that this section is informational only.

The addresses being used in the following examples are IPv4 addresses for private-use networks and for documentation as specified in [RFC6890]. However, the method described here is not limited to this NAT44 case.

The NAT binding table entries shown in the following examples do not include the flag indicating whether the restart procedure is supported or not. This flag is not relevant for these examples.

8.1. Single-homed Client to Single-homed Server

The internal client starts the association with the remote server via a four-way-handshake. Host A starts by sending a packet containing an INIT chunk.



```
INIT[Initiate-Tag = 1234]
10.0.0.1:1 -----> 203.0.113.1:2
    Rem-VTtag = 0
```

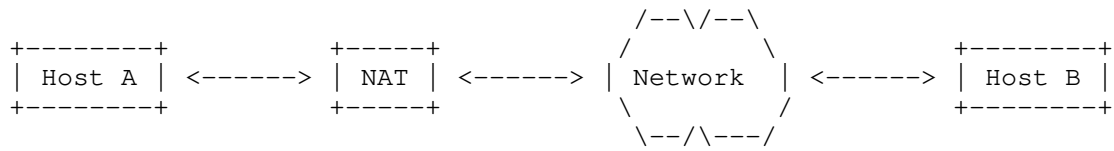
A NAT binding tabled entry is created, the source address is substituted and the packet is sent on:

NAT function creates entry:

NAT	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	0	2	10.0.0.1

```
INIT[Initiate-Tag = 1234]
192.0.2.1:1 -----> 203.0.113.1:2
    Rem-VTtag = 0
```

Host B receives the packet containing an INIT chunk and sends a packet containing an INIT ACK chunk with the NAT's Remote-address as destination address.



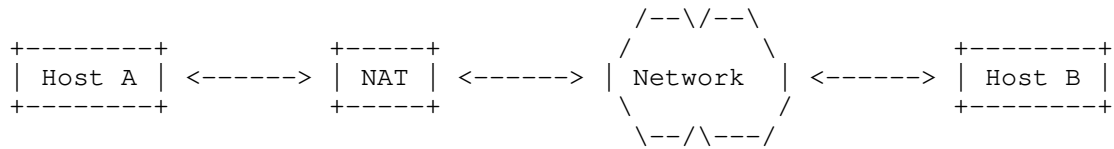
INIT ACK[Initiate-Tag = 5678]
 192.0.2.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

NAT function updates entry:

NAT	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	5678	2	10.0.0.1

INIT ACK[Initiate-Tag = 5678]
 10.0.0.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

The handshake finishes with a COOKIE ECHO acknowledged by a COOKIE ACK.



COOKIE ECHO
 10.0.0.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

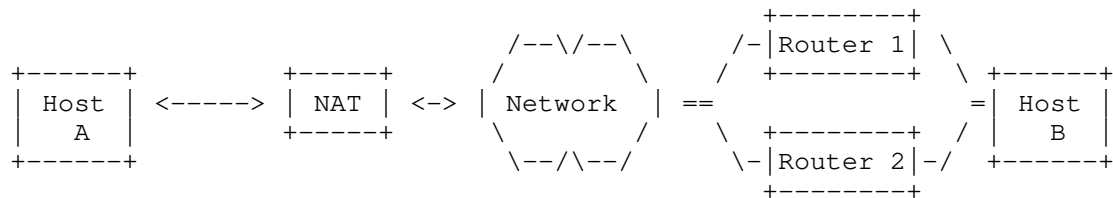
COOKIE ECHO
 192.0.2.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

COOKIE ACK
 192.0.2.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

COOKIE ACK
 10.0.0.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

8.2. Single-homed Client to Multi-homed Server

The internal client is single-homed whereas the remote server is multi-homed. The client (Host A) sends a packet containing an INIT chunk like in the single-homed case.



NAT	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
-----	-------------	-------------	-------------	-------------	-------------

```

INIT[Initiate-Tag = 1234]
10.0.0.1:1 ---> 203.0.113.1:2
Rem-VTag = 0
  
```

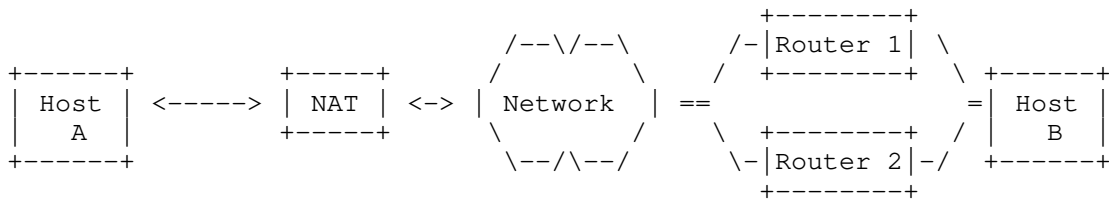
NAT function creates entry:

NAT	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	0	2	10.0.0.1

```

INIT[Initiate-Tag = 1234]
192.0.2.1:1 -----> 203.0.113.1:2
Rem-VTag = 0
  
```

The server (Host B) includes its two addresses in the INIT ACK chunk.



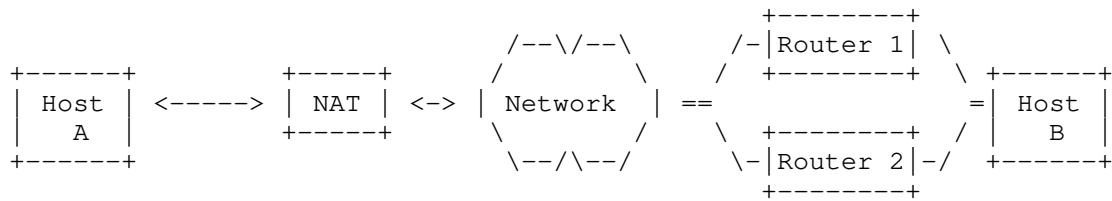
```
INIT ACK[Initiate-tag = 5678, IP-Addr = 203.0.113.129]
192.0.2.1:1 <----- 203.0.113.1:2
                        Int-VTag = 1234
```

The NAT function does not need to change the NAT binding table for the second address:

NAT	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	5678	2	10.0.0.1

```
INIT ACK[Initiate-Tag = 5678]
10.0.0.1:1 <--- 203.0.113.1:2
      Int-VTag = 1234
```

The handshake finishes with a COOKIE ECHO acknowledged by a COOKIE ACK.



COOKIE ECHO
 10.0.0.1:1 ---> 203.0.113.1:2
 Rem-VTag = 5678

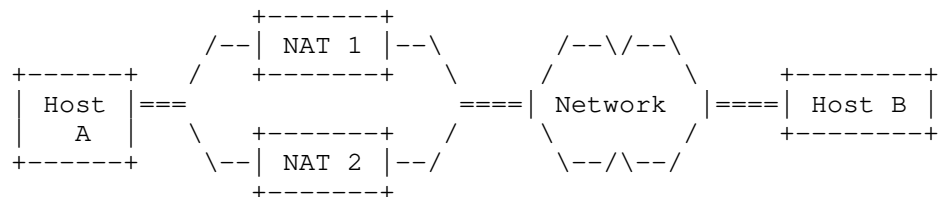
COOKIE ECHO
 192.0.2.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

COOKIE ACK
 192.0.2.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

COOKIE ACK
 10.0.0.1:1 <--- 203.0.113.1:2
 Int-VTag = 1234

8.3. Multihomed Client and Server

The client (Host A) sends a packet containing an INIT chunk to the server (Host B), but does not include the second address.



NAT 1	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr

INIT[Initiate-Tag = 1234]
 10.0.0.1:1 -----> 203.0.113.1:2
 Rem-VTag = 0

NAT function 1 creates entry:

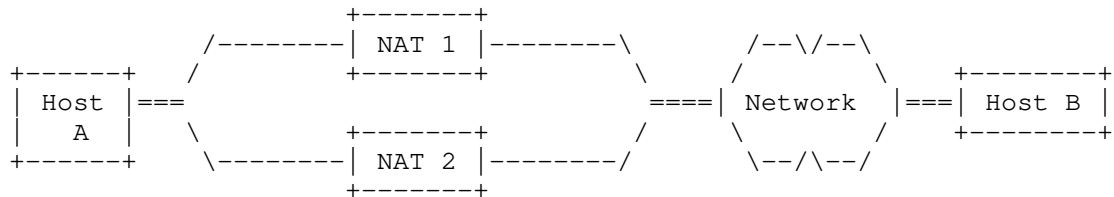
NAT 1	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	0	2	10.0.0.1

```

                                INIT[Initiate-Tag = 1234]
192.0.2.1:1 -----> 203.0.113.1:2
                                Rem-VTag = 0

```

Host B includes its second address in the INIT ACK.



```

INIT ACK[Initiate-Tag = 5678, IP-Addr = 203.0.113.129]
192.0.2.1:1 <----- 203.0.113.1:2
                                Int-VTag = 1234

```

NAT function 1 does not need to update the NAT binding table for the second address:

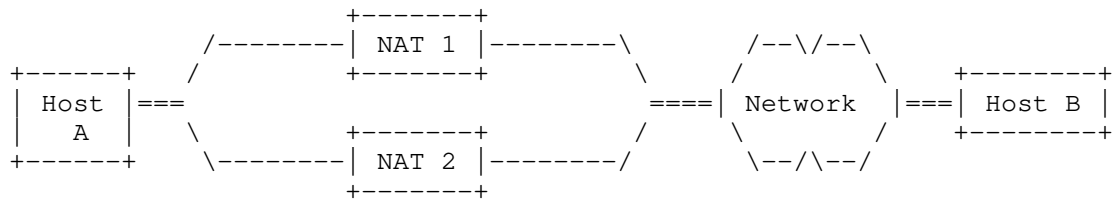
NAT 1	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	5678	2	10.0.0.1

```

INIT ACK[Initiate-Tag = 5678]
10.0.0.1:1 <----- 203.0.113.1:2
                                Int-VTag = 1234

```

The handshake finishes with a COOKIE ECHO acknowledged by a COOKIE ACK.



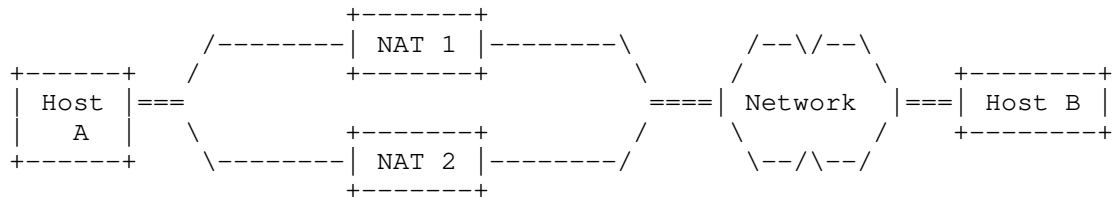
COOKIE ECHO
 10.0.0.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

COOKIE ECHO
 192.0.2.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

COOKIE ACK
 192.0.2.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

COOKIE ACK
 10.0.0.1:1 <----- 203.0.113.1:2
 Int-VTag = 1234

Host A announces its second address in an ASCONF chunk. The address parameter contains a wildcard address (0.0.0.0 or ::0) to indicate that the source address has to be added. The address parameter within the ASCONF chunk will also contain the pair of VTags (remote and internal) so that the NAT function can populate its NAT binding table entry completely with this single packet.



ASCONF [ADD-IP=0.0.0.0, INT-VTag=1234, Rem-VTag = 5678]
 10.1.0.1:1 -----> 203.0.113.129:2
 Rem-VTag = 5678

NAT function 2 creates a complete entry:

NAT 2	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	5678	2	10.1.0.1

```

ASCONF [ADD-IP, Int-VTag=1234, Rem-VTag = 5678]
192.0.2.129:1 -----> 203.0.113.129:2
                        Rem-VTag = 5678

```

```

                        ASCONF ACK
192.0.2.129:1 <----- 203.0.113.129:2
                        Int-VTag = 1234

```

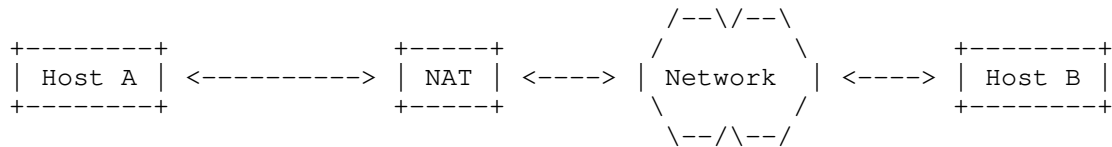
```

                        ASCONF ACK
10.1.0.1:1 <----- 203.0.113.129:2
                        Int-VTag = 1234

```

8.4. NAT Function Loses Its State

Association is already established between Host A and Host B, when the NAT function loses its state and obtains a new external address. Host A sends a DATA chunk to Host B.



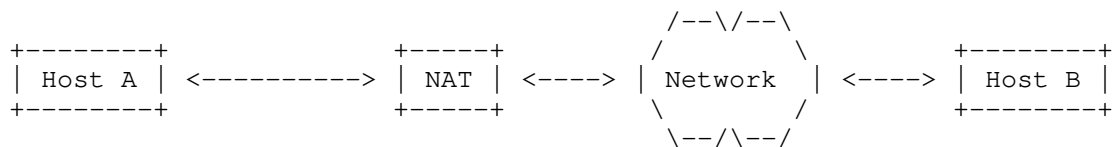
NAT	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr

```

                        DATA
10.0.0.1:1 -----> 203.0.113.1:2
                        Rem-VTag = 5678

```

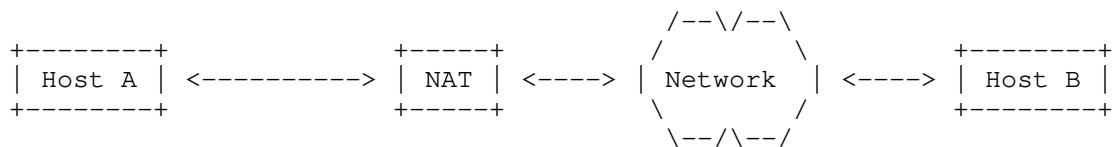
The NAT function cannot find an entry in the NAT binding table for the association. It sends a packet containing an ERROR chunk with the M bit set and the cause "NAT state missing".



```

ERROR [M bit, NAT state missing]
10.0.0.1:1 <----- 203.0.113.1:2
      Rem-VTag = 5678
  
```

On reception of the packet containing the ERROR chunk, Host A sends a packet containing an ASCONF chunk indicating that the former information has to be deleted and the source address of the actual packet added.



```

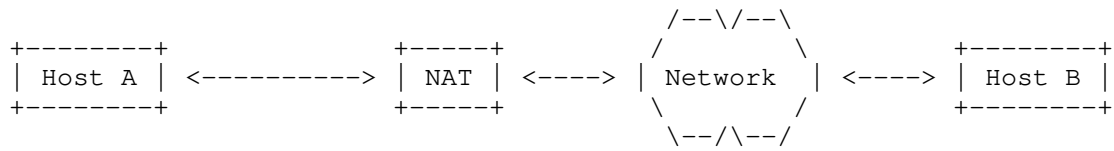
ASCONF [ADD-IP, DELETE-IP, Int-VTag=1234, Rem-VTag = 5678]
10.0.0.1:1 -----> 203.0.113.129:2
      Rem-VTag = 5678
  
```

NAT	+-----+ +-----+ +-----+ +-----+ +-----+				
	Int	Int	Rem	Rem	Int
	VTag	Port	VTag	Port	Addr
	+-----+ +-----+ +-----+ +-----+ +-----+				
	1234	1	5678	2	10.0.0.1
	+-----+ +-----+ +-----+ +-----+ +-----+				

```

ASCONF [ADD-IP, DELETE-IP, Int-VTag=1234, Rem-VTag = 5678]
      192.0.2.2:1 -----> 203.0.113.129:2
      Rem-VTag = 5678
  
```

Host B adds the new source address to this association and deletes all other addresses from this association.



ASCONF ACK
 192.0.2.2:1 <----- 203.0.113.129:2
 Int-VTag = 1234

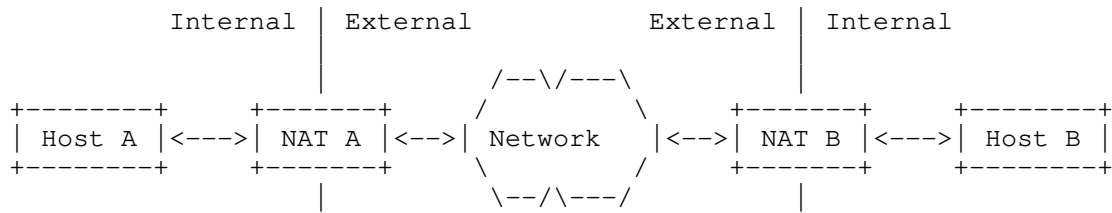
ASCONF ACK
 10.1.0.1:1 <----- 203.0.113.129:2
 Int-VTag = 1234

DATA
 10.0.0.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

DATA
 192.0.2.2:1 -----> 203.0.113.129:2
 Rem-VTag = 5678

8.5. Peer-to-Peer Communications

If two hosts, each of them behind a NAT function, want to communicate with each other, they have to get knowledge of the peer's external address. This can be achieved with a so-called rendezvous server. Afterwards the destination addresses are external, and the association is set up with the help of the INIT collision. The NAT functions create their entries according to their internal peer's point of view. Therefore, NAT function A's Internal-VTag and Internal-Port are NAT function B's Remote-VTag and Remote-Port, respectively. The naming (internal/remote) of the verification tag in the packet flow is done from the sending host's point of view.



NAT Binding Tables

NAT A	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
-------	-------------	-------------	-------------	-------------	-------------

NAT B	Int v-tag	Int port	Rem v-tag	Rem port	Int Addr
-------	--------------	-------------	--------------	-------------	-------------

```

INIT[Initiate-Tag = 1234]
10.0.0.1:1 --> 203.0.113.1:2
    Rem-VTag = 0
  
```

NAT function A creates entry:

NAT A	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	0	2	10.0.0.1

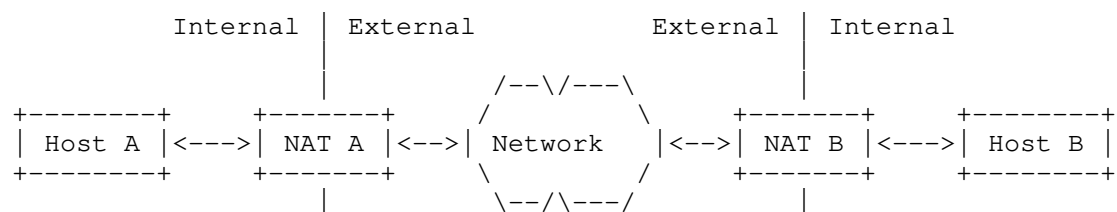
```

INIT[Initiate-Tag = 1234]
192.0.2.1:1 -----> 203.0.113.1:2
    Rem-VTag = 0
  
```

NAT function B processes the packet containing the INIT chunk, but cannot find an entry. The SCTP packet is silently discarded and leaves the NAT binding table of NAT function B unchanged.

NAT B	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
-------	-------------	-------------	-------------	-------------	-------------

Now Host B sends a packet containing an INIT chunk, which is processed by NAT function B. Its parameters are used to create an entry.



```

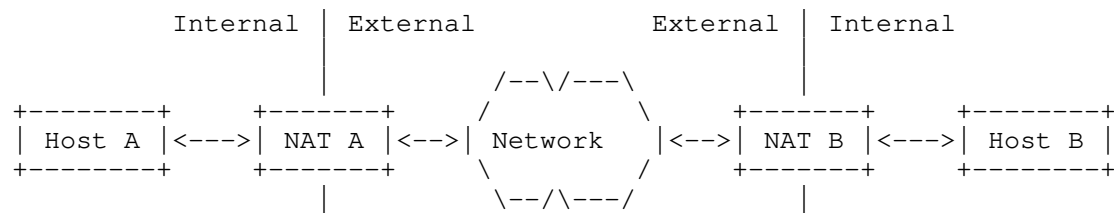
INIT[Initiate-Tag = 5678]
192.0.2.1:1 <-- 10.1.0.1:2
Rem-VTag = 0
  
```

NAT B	Int	Int	Rem	Rem	Int
	VTag	Port	VTag	Port	Addr
	5678	2	0	1	10.1.0.1

```

INIT[Initiate-Tag = 5678]
192.0.2.1:1 <----- 203.0.113.1:2
Rem-VTag = 0
  
```

NAT function A processes the packet containing the INIT chunk. As the outgoing packet containing an INIT chunk of Host A has already created an entry, the entry is found and updated:

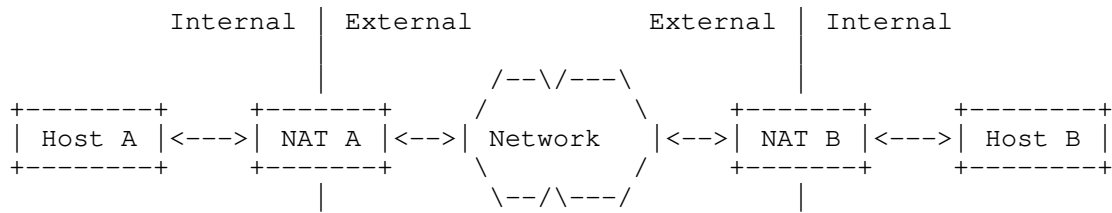


VTag != Int-VTag, but Rem-VTag == 0, find entry.

NAT A	Int VTag	Int Port	Rem VTag	Rem Port	Int Addr
	1234	1	5678	2	10.0.0.1

```
INIT[Initiate-tag = 5678]
10.0.0.1:1 <-- 203.0.113.1:2
    Rem-VTag = 0
```

Host A sends a packet containing an INIT ACK chunk, which can pass through NAT function B:



```

INIT ACK[Initiate-Tag = 1234]
10.0.0.1:1 --> 203.0.113.1:2
    Rem-VTag = 5678

```

```

        INIT ACK[Initiate-Tag = 1234]
192.0.2.1:1 -----> 203.0.113.1:2
        Rem-VTag = 5678

```

NAT function B updates entry:

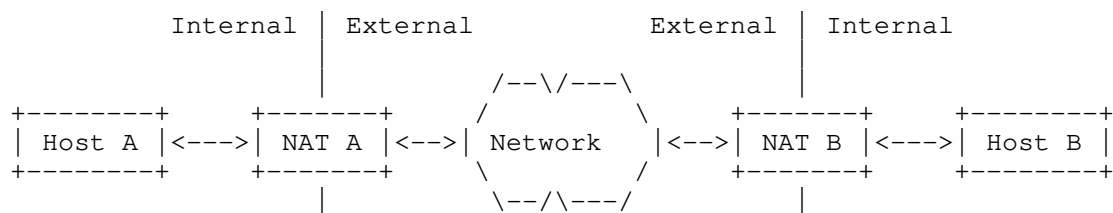
NAT B	Int	Int	Rem	Rem	Int
	VTag	Port	VTag	Port	Addr
	5678	2	1234	1	10.1.0.1

```

INIT ACK[Initiate-Tag = 1234]
192.0.2.1:1 --> 10.1.0.1:2
    Rem-VTag = 5678

```

The lookup for COOKIE ECHO and COOKIE ACK is successful.



COOKIE ECHO
 192.0.2.1:1 <-- 10.1.0.1:2
 Rem-VTag = 1234

COOKIE ECHO
 192.0.2.1:1 <----- 203.0.113.1:2
 Rem-VTag = 1234

COOKIE ECHO
 10.0.0.1:1 <-- 203.0.113.1:2
 Rem-VTag = 1234

COOKIE ACK
 10.0.0.1:1 --> 203.0.113.1:2
 Rem-VTag = 5678

COOKIE ACK
 192.0.2.1:1 -----> 203.0.113.1:2
 Rem-VTag = 5678

COOKIE ACK
 192.0.2.1:1 --> 10.1.0.1:2
 Rem-VTag = 5678

9. Socket API Considerations

This section describes how the socket API defined in [RFC6458] is extended to provide a way for the application to control NAT friendliness.

Please note that this section is informational only.

A socket API implementation based on [RFC6458] is extended by supporting one new read/write socket option.

9.1. Get or Set the NAT Friendliness (SCTP_NAT_FRIENDLY)

This socket option uses the option_level IPPROTO_SCTP and the option_name SCTP_NAT_FRIENDLY. It can be used to enable/disable the NAT friendliness for future associations and retrieve the value for future and specific ones.

```
struct sctp_assoc_value {  
    sctp_assoc_t assoc_id;  
    uint32_t assoc_value;  
};
```

assoc_id

This parameter is ignored for one-to-one style sockets. For one-to-many style sockets the application can fill in an association identifier or SCTP_FUTURE_ASSOC for this query. It is an error to use SCTP_{CURRENT|ALL}_ASSOC in assoc_id.

assoc_value

A non-zero value indicates a NAT-friendly mode.

10. IANA Considerations

[NOTE to RFC-Editor: "RFCXXXX" is to be replaced by the RFC number you assign this document.]

[NOTE to RFC-Editor: The requested values for the chunk type and the chunk parameter types are tentative and to be confirmed by IANA.]

This document (RFCXXXX) is the reference for all registrations described in this section. The requested changes are described below.

10.1. New Chunk Flags for Two Existing Chunk Types

As defined in [RFC6096] two chunk flags have to be assigned by IANA for the ERROR chunk. The requested value for the T bit is 0x01 and for the M bit is 0x02.

This requires an update of the "ERROR Chunk Flags" registry for SCTP:

ERROR Chunk Flags

Chunk Flag Value	Chunk Flag Name	Reference
0x01	T bit	[RFCXXXX]
0x02	M bit	[RFCXXXX]
0x04	Unassigned	
0x08	Unassigned	
0x10	Unassigned	
0x20	Unassigned	
0x40	Unassigned	
0x80	Unassigned	

Table 2

As defined in [RFC6096] one chunk flag has to be assigned by IANA for the ABORT chunk. The requested value of the M bit is 0x02.

This requires an update of the "ABORT Chunk Flags" registry for SCTP:

ABORT Chunk Flags

Chunk Flag Value	Chunk Flag Name	Reference
0x01	T bit	[RFC4960]
0x02	M bit	[RFCXXXX]
0x04	Unassigned	
0x08	Unassigned	
0x10	Unassigned	
0x20	Unassigned	
0x40	Unassigned	
0x80	Unassigned	

Table 3

10.2. Three New Error Causes

Three error causes have to be assigned by IANA. It is requested to use the values given below.

This requires three additional lines in the "Error Cause Codes" registry for SCTP:

Error Cause Codes

Value	Cause Code	Reference
176	VTag and Port Number Collision	[RFCXXXX]
177	Missing State	[RFCXXXX]
178	Port Number Collision	[RFCXXXX]

Table 4

10.3. Two New Chunk Parameter Types

Two chunk parameter types have to be assigned by IANA. IANA is requested to assign these values from the pool of parameters with the upper two bits set to '11' and to use the values given below.

This requires two additional lines in the "Chunk Parameter Types" registry for SCTP:

Chunk Parameter Types

ID Value	Chunk Parameter Type	Reference
49159	Disable Restart (0xC007)	[RFCXXXX]
49160	VTags (0xC008)	[RFCXXXX]

Table 5

10.4. One New URI

An URI in the "ns" subregistry within the "IETF XML" registry has to be assigned by IANA ([RFC3688]):

URI: urn:ietf:params:xml:ns:yang:ietf-nat-sctp
 Registrant Contact: The IESG.
 XML: N/A; the requested URI is an XML namespace.

10.5. One New YANG Module

An YANG module in the "YANG Module Names" subregistry within the "YANG Parameters" registry has to be assigned by IANA ([RFC6020]):

Name: ietf-nat-sctp
 Namespace: urn:ietf:params:xml:ns:yang:ietf-nat-sctp
 Maintained by IANA: N
 Prefix: nat-sctp
 Reference: RFCXXXX

11. Security Considerations

State maintenance within a NAT function is always a subject of possible Denial Of Service attacks. This document recommends that at a minimum a NAT function runs a timer on any SCTP state so that old association state can be cleaned up.

Generic issues related to address sharing are discussed in [RFC6269] and apply to SCTP as well.

For SCTP endpoints not disabling the restart procedure, this document does not add any additional security considerations to the ones given in [RFC4960], [RFC4895], and [RFC5061].

SCTP endpoints disabling the restart procedure, need to monitor the status of all associations to mitigate resource exhaustion attacks by establishing a lot of associations sharing the same IP addresses and port numbers.

In any case, SCTP is protected by the verification tags and the usage of [RFC4895] against off-path attackers.

For IP-level fragmentation and reassembly related issues see [RFC4963].

The YANG module specified in this document defines a schema for data that is designed to be accessed via network management protocols such as NETCONF [RFC6241] or RESTCONF [RFC8040]. The lowest NETCONF layer is the secure transport layer, and the mandatory-to-implement secure transport is Secure Shell (SSH) [RFC6242]. The lowest RESTCONF layer is HTTPS, and the mandatory-to-implement secure transport is TLS [RFC8446].

The Network Configuration Access Control Model (NACM) [RFC8341] provides the means to restrict access for particular NETCONF or RESTCONF users to a preconfigured subset of all available NETCONF or RESTCONF protocol operations and content.

All data nodes defined in the YANG module that can be created, modified, and deleted (i.e., config true, which is the default) are considered sensitive. Write operations (e.g., edit-config) applied to these data nodes without proper protection can negatively affect network operations. An attacker who is able to access the SCTP NAT function can undertake various attacks, such as:

- * Setting a low timeout for SCTP mapping entries to cause failures to deliver incoming SCTP packets.
- * Instantiating mapping entries to cause NAT collision.

12. Normative References

- [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>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC4895] Tuexen, M., Stewart, R., Lei, P., and E. Rescorla, "Authenticated Chunks for the Stream Control Transmission Protocol (SCTP)", RFC 4895, DOI 10.17487/RFC4895, August 2007, <<https://www.rfc-editor.org/info/rfc4895>>.
- [RFC4960] Stewart, R., Ed., "Stream Control Transmission Protocol", RFC 4960, DOI 10.17487/RFC4960, September 2007, <<https://www.rfc-editor.org/info/rfc4960>>.
- [RFC5061] Stewart, R., Xie, Q., Tuexen, M., Maruyama, S., and M. Kozuka, "Stream Control Transmission Protocol (SCTP) Dynamic Address Reconfiguration", RFC 5061, DOI 10.17487/RFC5061, September 2007, <<https://www.rfc-editor.org/info/rfc5061>>.
- [RFC6020] Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.
- [RFC6096] Tuexen, M. and R. Stewart, "Stream Control Transmission Protocol (SCTP) Chunk Flags Registration", RFC 6096, DOI 10.17487/RFC6096, January 2011, <<https://www.rfc-editor.org/info/rfc6096>>.
- [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>>.
- [RFC8512] Boucadair, M., Ed., Sivakumar, S., Jacquenet, C., Vinapamula, S., and Q. Wu, "A YANG Module for Network Address Translation (NAT) and Network Prefix Translation (NPT)", RFC 8512, DOI 10.17487/RFC8512, January 2019, <<https://www.rfc-editor.org/info/rfc8512>>.

13. Informative References

- [DOI_10.1145_1496091.1496095]
Hayes, D., But, J., and G. Armitage, "Issues with network address translation for SCTP", ACM SIGCOMM Computer Communication Review Vol. 39, pp. 23-33, DOI 10.1145/1496091.1496095, December 2008, <<https://doi.org/10.1145/1496091.1496095>>.
- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, DOI 10.17487/RFC0793, September 1981, <<https://www.rfc-editor.org/info/rfc793>>.
- [RFC3022] Srisuresh, P. and K. Egevang, "Traditional IP Network Address Translator (Traditional NAT)", RFC 3022, DOI 10.17487/RFC3022, January 2001, <<https://www.rfc-editor.org/info/rfc3022>>.
- [RFC4787] Audet, F., Ed. and C. Jennings, "Network Address Translation (NAT) Behavioral Requirements for Unicast UDP", BCP 127, RFC 4787, DOI 10.17487/RFC4787, January 2007, <<https://www.rfc-editor.org/info/rfc4787>>.
- [RFC4963] Heffner, J., Mathis, M., and B. Chandler, "IPv4 Reassembly Errors at High Data Rates", RFC 4963, DOI 10.17487/RFC4963, July 2007, <<https://www.rfc-editor.org/info/rfc4963>>.
- [RFC5382] Guha, S., Ed., Biswas, K., Ford, B., Sivakumar, S., and P. Srisuresh, "NAT Behavioral Requirements for TCP", BCP 142, RFC 5382, DOI 10.17487/RFC5382, October 2008, <<https://www.rfc-editor.org/info/rfc5382>>.
- [RFC5508] Srisuresh, P., Ford, B., Sivakumar, S., and S. Guha, "NAT Behavioral Requirements for ICMP", BCP 148, RFC 5508, DOI 10.17487/RFC5508, April 2009, <<https://www.rfc-editor.org/info/rfc5508>>.
- [RFC6056] Larsen, M. and F. Gont, "Recommendations for Transport-Protocol Port Randomization", BCP 156, RFC 6056, DOI 10.17487/RFC6056, January 2011, <<https://www.rfc-editor.org/info/rfc6056>>.
- [RFC6146] Bagnulo, M., Matthews, P., and I. van Beijnum, "Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers", RFC 6146, DOI 10.17487/RFC6146, April 2011, <<https://www.rfc-editor.org/info/rfc6146>>.

- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6242] Wasserman, M., "Using the NETCONF Protocol over Secure Shell (SSH)", RFC 6242, DOI 10.17487/RFC6242, June 2011, <<https://www.rfc-editor.org/info/rfc6242>>.
- [RFC6269] Ford, M., Ed., Boucadair, M., Durand, A., Levis, P., and P. Roberts, "Issues with IP Address Sharing", RFC 6269, DOI 10.17487/RFC6269, June 2011, <<https://www.rfc-editor.org/info/rfc6269>>.
- [RFC6333] Durand, A., Droms, R., Woodyatt, J., and Y. Lee, "Dual-Stack Lite Broadband Deployments Following IPv4 Exhaustion", RFC 6333, DOI 10.17487/RFC6333, August 2011, <<https://www.rfc-editor.org/info/rfc6333>>.
- [RFC6458] Stewart, R., Tuexen, M., Poon, K., Lei, P., and V. Yasevich, "Sockets API Extensions for the Stream Control Transmission Protocol (SCTP)", RFC 6458, DOI 10.17487/RFC6458, December 2011, <<https://www.rfc-editor.org/info/rfc6458>>.
- [RFC6890] Cotton, M., Vegoda, L., Bonica, R., Ed., and B. Haberman, "Special-Purpose IP Address Registries", BCP 153, RFC 6890, DOI 10.17487/RFC6890, April 2013, <<https://www.rfc-editor.org/info/rfc6890>>.
- [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>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC7857] Penno, R., Perreault, S., Boucadair, M., Ed., Sivakumar, S., and K. Naito, "Updates to Network Address Translation (NAT) Behavioral Requirements", BCP 127, RFC 7857, DOI 10.17487/RFC7857, April 2016, <<https://www.rfc-editor.org/info/rfc7857>>.

- [RFC8040] Bierman, A., Bjorklund, M., and K. Watsen, "RESTCONF Protocol", RFC 8040, DOI 10.17487/RFC8040, January 2017, <<https://www.rfc-editor.org/info/rfc8040>>.
- [RFC8340] Bjorklund, M. and L. Berger, Ed., "YANG Tree Diagrams", BCP 215, RFC 8340, DOI 10.17487/RFC8340, March 2018, <<https://www.rfc-editor.org/info/rfc8340>>.
- [RFC8341] Bierman, A. and M. Bjorklund, "Network Configuration Access Control Model", STD 91, RFC 8341, DOI 10.17487/RFC8341, March 2018, <<https://www.rfc-editor.org/info/rfc8341>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/info/rfc8446>>.
- [RFC8900] Bonica, R., Baker, F., Huston, G., Hinden, R., Troan, O., and F. Gont, "IP Fragmentation Considered Fragile", BCP 230, RFC 8900, DOI 10.17487/RFC8900, September 2020, <<https://www.rfc-editor.org/info/rfc8900>>.

Acknowledgments

The authors wish to thank Mohamed Boucadair, Gorrry Fairhurst, Bryan Ford, David Hayes, Alfred Hines, Karen E. E. Nielsen, Henning Peters, Maksim Proshin, Timo Völker, Dan Wing, and Qiaobing Xie for their invaluable comments.

In addition, the authors wish to thank David Hayes, Jason But, and Grenville Armitage, the authors of [DOI_10.1145_1496091.1496095], for their suggestions.

The authors also wish to thank Mohamed Boucadair for contributing the text related to the YANG module.

Authors' Addresses

Randall R. Stewart
Netflix, Inc.
Chapin, SC 29036
United States of America

Email: randall@lakerest.net

Michael Tüxen
Münster University of Applied Sciences
Stegerwaldstrasse 39
48565 Steinfurt
Germany

Email: tuexen@fh-muenster.de

Irene Rüngeler
Münster University of Applied Sciences
Stegerwaldstrasse 39
48565 Steinfurt
Germany

Email: i.ruengeler@fh-muenster.de

TSVWG
Internet Draft
Intended status: Best Current Practice
Expires: October 2015

J. Touch
USC/ISI
April 24, 2015

Recommendations on Using Assigned Transport Port Numbers
draft-ietf-tsvwg-port-use-11.txt

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), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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

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

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

This Internet-Draft will expire on October 24, 2015.

Copyright Notice

Copyright (c) 2015 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Abstract

This document provides recommendations to application and service protocol designers on how to use the assigned transport protocol port number space and when to request a port assignment from IANA. It provides designer guidelines on how to interact with the IANA processes defined in RFC6335, thus serving to complement (but not update) that document.

Table of Contents

1. Introduction.....	2
2. Conventions used in this document.....	3
3. History.....	3
4. Current Port Number Use.....	5
5. What is a Port Number?.....	5
6. Conservation.....	7
6.1. Guiding Principles.....	7
6.2. Firewall and NAT Considerations.....	8
7. Considerations for Requesting Port Number Assignments.....	9
7.1. Is a port number assignment necessary?.....	9
7.2. How Many Assigned Port Numbers?.....	11
7.3. Picking an Assigned Port Number.....	12
7.4. Support for Security.....	13
7.5. Support for Future Versions.....	14
7.6. Transport Protocols.....	15
7.7. When to Request an Assignment.....	16
7.8. Squatting.....	17
7.9. Other Considerations.....	18
8. Security Considerations.....	18
9. IANA Considerations.....	19
10. References.....	19
10.1. Normative References.....	19
10.2. Informative References.....	20
11. Acknowledgments.....	22

1. Introduction

This document provides information and advice to application and service designers on the use of assigned transport port numbers. It provides a detailed historical background of the evolution of transport port numbers and their multiple meanings. It also provides specific recommendations to designers on how to use assigned port numbers. Note that this document provides information to potential port number applicants that complements the IANA process described in BCP165 [RFC6335], but it does not change any of the port number

assignment procedures described therein. This document is intended to address concerns typically raised during Expert Review of assigned port number applications, but it is not intended to bind those reviews. RFC 6335 also describes the interaction between port experts and port requests in IETF consensus document. Authors of IETF consensus documents should nevertheless follow the advice in this document and can expect comment on their port requests from the port experts during IETF last call or at other times when review is explicitly sought.

2. Conventions used in this document

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC-2119 [RFC2119].

In this document, these words will appear with that interpretation only when in ALL CAPS. Lower case uses of these words are not to be interpreted as carrying RFC-2119 significance.

In this document, the characters ">>" preceding an indented line(s) indicates a statement using the key words listed above. This convention aids reviewers in quickly identifying or finding requirements for registration and recommendations for use of port numbers in this RFC.

3. History

The term 'port' was first used in [RFC33] to indicate a simplex communication path from an individual process and originally applied to only the Network Control Program (NCP) connection-oriented protocol. At a meeting described in [RFC37], an idea was presented to decouple connections between processes and links that they use as paths, and thus to include numeric source and destination socket identifiers in packets. [RFC38] provides further detail, describing how processes might have more than one of these paths and that more than one path may be active at a time. As a result, there was the need to add a process identifier to the header of each message so that incoming messages could be demultiplexed to the appropriate process. [RFC38] further suggested that 32 bit numbers would be used for these identifiers. [RFC48] discusses the current notion of listening on a specific port number, but does not discuss the issue of port number determination. [RFC61] notes that the challenge of knowing the appropriate port numbers is "left to the processes" in general, but introduces the concept of a "well-known" port number for common services.

[RFC76] proposed a "telephone book" by which an index would allow port numbers to be used by name, but still assumed that both source and destination port numbers are fixed by such a system. [RFC333] proposed that a port number pair, rather than an individual port number, would be used on both sides of the connection for demultiplexing messages. This is the final view in [RFC793] (and its predecessors, including [IEN112]), and brings us to their current meaning. [RFC739] introduced the notion of generic reserved port numbers for groups of protocols, such as "any private RJE server" [RFC739]. Although the overall range of such port numbers was (and remains) 16 bits, only the first 256 (high 8 bits cleared) in the range were considered assigned.

[RFC758] is the first to describe port numbers as being used for TCP (previous RFCs all refer to only NCP). It includes a list of such well-known port numbers, as well as describing ranges used for different purposes:

Decimal	Octal	
---------	-------	--

0-63	0-77	Network Wide Standard Function
64-127	100-177	Hosts Specific Functions
128-223	200-337	Reserved for Future Use
224-255	340-377	Any Experimental Function

In [RFC820] those range meanings disappeared, and a single list of number assignments is presented. This is also the first time that port numbers are described as applying to a connectionless transport (UDP) rather than only connection-oriented transports.

By [RFC900] the ranges appeared as decimal numbers rather than the octal ranges used previously. [RFC1340] increased this range from 0..255 to 0..1023, and began to list TCP and UDP port number assignments individually (although the assumption was that once assigned a port number applies to all transport protocols, including TCP, UDP, recently SCTP and DCCP, as well as ISO-TP4 for a brief period in the early 1990s). [RFC1340] also established the Registered range of 1024-59151, though it notes that it is not controlled by the IANA at that point. The list provided by [RFC1700] in 1994 remained the standard until it was declared replaced by an on-line version, as of [RFC3232] in 2002.

4. Current Port Number Use

RFC6335 indicates three ranges of port number assignments:

Binary	Hex	

0-1023	0x0000-0x03FF	System (also Well-Known)
1024-49151	0x0400-0xBFFF	User (also Registered)
49152-65535	0xC000-0xFFFF	Dynamic (also Private)

System (also Well-Known) encompasses the range 0..1023. On some systems, use of these port numbers requires privileged access, e.g., that the process run as 'root' (i.e., as a privileged user), which is why these are referred to as System port numbers. The port numbers from 1024..49151 denotes non-privileged services, known as User (also Registered), because these port numbers do not run with special privileges. Dynamic (also Private) port numbers are not assigned.

Both System and User port numbers are assigned through IANA, so both are sometimes called 'registered port numbers'. As a result, the term 'registered' is ambiguous, referring either to the entire range 0-49151 or to the User port numbers. Complicating matters further, System port numbers do not always require special (i.e., 'root') privilege. For clarity, the remainder of this document refers to the port number ranges as System, User, and Dynamic, to be consistent with IANA process [RFC6335].

5. What is a Port Number?

A port number is a 16-bit number used for two distinct purposes:

- o Demultiplexing transport endpoint associations within an end host
- o Identifying a service

The first purpose requires that each transport endpoint association (e.g., TCP connection or UDP pairwise association) using a given transport between a given pair of IP addresses use a different pair of port numbers, but does not require either coordination or registration of port number use. It is the second purpose that drives the need for a common registry.

Consider a user wanting to run a web server. That service could run on any port number, provided that all clients knew what port number to use to access that service at that host. Such information can be explicitly distributed - for example, by putting it in the URI:

`http://www.example.com:51509/`

Ultimately, the correlation of a service with a port number is an agreement between just the two endpoints of the association. A web server can run on port number 53, which might appear as DNS traffic to others but will connect to browsers that know to use port number 53 rather than 80.

As a concept, a service is the combination of ISO Layers 5-7 that represents an application protocol capability. For example www (port number 80) is a service that uses HTTP as an application protocol and provides access to a web server [RFC7230]. However, it is possible to use HTTP for other purposes, such as command and control. This is why some current services (HTTP, e.g.) are a bit overloaded - they describe not only the application protocol, but a particular service.

IANA assigns port numbers so that Internet endpoints do not need pairwise, explicit coordination of the meaning of their port numbers. This is the primary reason for requesting port number assignment by IANA - to have a common agreement between all endpoints on the Internet as to the default meaning of a port number, which provides the endpoints with a default port number for a particular protocol or service.

Port numbers are sometimes used by intermediate devices on a network path, either to monitor available services, to monitor traffic (e.g., to indicate the data contents), or to intercept traffic (to block, proxy, relay, aggregate, or otherwise process it). In each case, the intermediate device interprets traffic based on the port number. It is important to recognize that any interpretation of port numbers - except at the endpoints - may be incorrect, because port numbers are meaningful only at the endpoints. Further, port numbers may not be visible to these intermediate devices, such as when the transport protocol is encrypted (as in network- or link-layer tunnels), or when a packet is fragmented (in which case only the first fragment has the port number information). Such port number invisibility may interfere with these in-network port number-based capabilities.

Port numbers can also be used for other purposes. Assigned port numbers can simplify end system configuration, so that individual

installations do not need to coordinate their use of arbitrary port numbers. Such assignments may also have the effect of simplifying firewall management, so that a single, fixed firewall configuration can either permit or deny a service that uses the assigned ports.

It is useful to differentiate a port number from a service name. The former is a numeric value that is used directly in transport protocol headers as a demultiplexing and service identifier. The latter is primarily a user convenience, where the default map between the two is considered static and resolved using a cached index. This document focuses on the former because it is the fundamental network resource. Dynamic maps between the two, i.e., using DNS SRV records, are discussed further in Section 7.1.

6. Conservation

Assigned port numbers are a limited resource that is globally shared by the entire Internet community. As of 2014, approximately 5850 TCP and 5570 UDP port numbers have been assigned out of a total range of 49151. As a result of past conservation, current assigned port use is small and the current rate of assignment avoids the need for transition to larger number spaces. This conservation also helps avoid the need for IANA to rely on assigned port number reclamation, which is practically impossible even though procedurally permitted [RFC6335].

IANA aims to assign only one port number per service, including variants [RFC6335], but there are other benefits to using fewer port numbers for a given service. Use of multiple assigned port numbers can make applications more fragile, especially when firewalls block a subset of those port numbers or use ports numbers to route or prioritize traffic differently. As a result:

>> Each assigned port requested MUST be justified by the applicant as an independently useful service.

6.1. Guiding Principles

This document provides recommendations for users that also help conserve assigned port number space. Again, this document does not update BCP165 [RFC6335], which describes the IANA procedures for managing assigned transport port numbers and services. Assigned port number conservation is based on a number of basic principles:

- o A single assigned port number can support different functions over separate endpoint associations, determined using in-band information. An FTP data connection can transfer binary or text files, the latter translating line-terminators, as indicated in-band over the control port number [RFC959].
- o A single assigned port number can indicate the Dynamic port number(s) on which different capabilities are supported, as with passive-mode FTP [RFC959].
- o Several existing services can indicate the Dynamic port number(s) on which other services are supported, such as with mDNS and portmapper [RFC1833] [RFC6762] [RFC6763].
- o Copies of some existing services can be differentiated using in-band information (e.g., URIs in HTTP Host field and TLS Server Name Indication extension) [RFC7230] [RFC6066].
- o Services requiring varying performance properties can already be supported using separate endpoint associations (connections or other associations), each configured to support the desired properties. E.g., a high-speed and low-speed variant can be determined within the service using the same assigned port.

Assigned port numbers are intended to differentiate services, not variations of performance, replicas, pairwise endpoint associations, or payload types. Assigned port numbers are also a small space compared to other Internet number spaces; it is never appropriate to consume assigned port numbers to conserve larger spaces such as IP addresses, especially where copies of a service represent different endpoints.

6.2. Firewall and NAT Considerations

Ultimately, port numbers indicate services only to the endpoints, and any intermediate device that assigns meaning to a value can be incorrect. End systems might agree to run web services (HTTP) over port number 53 (typically used for DNS) rather than port number 80, at which point a firewall that blocks port number 80 but permits port number 53 would not have the desired effect. Nonetheless, assigned port numbers are often used to help configure firewalls and other port-based systems for access control.

Using Dynamic port numbers, or explicitly-indicated port numbers indicated in-band over another service (such as with FTP) often complicates firewall and NAT interactions [RFC959]. FTP over firewalls often requires direct support for deep-packet inspection

(to snoop for the Dynamic port number for the NAT to correctly map) or passive-mode FTP (in which both connections are opened from the client side).

7. Considerations for Requesting Port Number Assignments

Port numbers are assigned by IANA by a set of documented procedures [RFC6335]. The following section describes the steps users can take to help assist with responsible use of assigned port numbers, and with preparing an application for a port number assignment.

7.1. Is a port number assignment necessary?

First, it is useful to consider whether a port number assignment is required. In many cases, a new number assignment may not be needed, for example:

- o Is this really a new service, or can an existing service suffice?
- o Is this an experimental service [RFC3692]? If so, consider using the current experimental ports [RFC2780].
- o Is this service independently useful? Some systems are composed from collections of different service capabilities, but not all component functions are useful as independent services. Port numbers are typically shared among the smallest independently-useful set of functions. Different service uses or properties can be supported in separate pairwise endpoint associations after an initial negotiation, e.g., to support software decomposition.
- o Can this service use a Dynamic port number that is coordinated out-of-band, e.g.:
 - o By explicit configuration of both endpoints.
 - o By internal mechanisms within the same host (e.g., a configuration file, indicated within a URI, or using interprocess communication).
- o Using information exchanged on a related service: FTP, SIP, etc. [RFC959] [RFC3261].
- o Using an existing port discovery service: portmapper, mDNS, etc. [RFC1833] [RFC6762] [RFC6763].

There are a few good examples of reasons that more directly suggest that not only is a port number assignment not necessary, but it is directly counter-indicated:

- o Assigned port numbers are not intended to differentiate performance variations within the same service, e.g., high-speed vs. ordinary speed. Performance variations can be supported within a single assigned port number in context of separate pairwise endpoint associations.
- o Additional assigned port numbers are not intended to replicate an existing service. For example, if a device is configured to use a typical web browser then it the port number used for that service is a copy of the http service that is already assigned to port number 80 and does not warrant a new assignment. However, an automated system that happens to use HTTP framing - but is not primarily accessed by a browser - might be a new service. A good way to tell is "can an unmodified client of the existing service interact with the proposed service"? If so, that service would be a copy of an existing service and would not merit a new assignment.
- o Assigned port numbers not intended for intra-machine communication. Such communication can already be supported by internal mechanisms (interprocess communication, shared memory, shared files, etc.). When Internet communication within a host is desired, the server can bind to a Dynamic port that is indicated to the client using these internal mechanisms.
- o Separate assigned port numbers are not intended for insecure versions of existing (or new) secure services. A service that already requires security would be made more vulnerable by having the same capability accessible without security.

Note that the converse is different, i.e., it can be useful to create a new, secure service that replicates an existing insecure service on a new port number assignment. This can be necessary when the existing service is not backward-compatible with security enhancements, such as the use of TLS [RFC5246] or DTLS [RFC6347].

- o Assigned port numbers are not intended for indicating different service versions. Version differentiation should be handled in-band, e.g., using a version number at the beginning of an association (e.g., connection or other transaction). This may not be possible with legacy assignments, but all new services should incorporate support for version indication.

Some services may not need assigned port numbers at all, e.g., SIP allows voice calls to use Dynamic ports [RFC3261]. Some systems can register services in the DNS, using SRV entries. These services can be discovered by a variety of means, including mDNS, or via direct query [RFC6762] [RFC6763]. In such cases, users can more easily request a SRV name, which are assigned first-come, first-served from a much larger namespace.

IANA assigns port numbers, but this assignment is typically used only for servers, i.e., the host that listens for incoming connections or other associations. Clients, i.e., hosts that initiate connections or other associations, typically refer to those assigned port numbers but do not need port number assignments for their endpoint.

Finally, an assigned port number is not a guarantee of exclusive use. Traffic for any service might appear on any port number, due to misconfiguration or deliberate misuse. Application and service designers are encouraged to validate traffic based on its content.

7.2. How Many Assigned Port Numbers?

As noted earlier, systems might require a single port number assignment, but rarely require multiple port numbers. There are a variety of known ways to reduce assigned port number consumption. Although some may be cumbersome or inefficient, they are nearly always preferable to consuming additional port number assignments.

Such techniques include:

- o Use of a discovery service, either a shared service (mDNS), or a discovery service for a given system [RFC6762] [RFC6763].
- o Multiplex packet types using in-band information, either on a per-message or per-connection basis. Such demultiplexing can even hand-off different messages and connections among different processes, such as is done with FTP [RFC959].

There are some cases where NAT and firewall traversal are significantly improved by having an assigned port number. Although

NAT traversal protocols supporting automatic configuration have been proposed and developed (e.g., STUN [RFC5389], TURN [RFC5766], and ICE [RFC5245]), not all application and service designers can rely on their presence as of yet.

In the past, some services were assigned multiple port numbers or sometimes fairly large port ranges (e.g., X11). This occurred for a variety of reasons: port number conservation was not as widely appreciated, assignments were not as ardently reviewed, etc. This no longer reflects current practice and such assignments are not considered to constitute a precedent for future assignments.

7.3. Picking an Assigned Port Number

Given a demonstrated need for a port number assignment, the next question is how to pick the desired port number. An application for a port number assignment does not need to include a desired port number; in that case, IANA will select from those currently available.

Users should consider whether the requested port number is important. For example, would an assignment be acceptable if IANA picked the port number value? Would a TCP (or other transport protocol) port number assignment be useful by itself? If so, a port number can be assigned to a service for one transport protocol where it is already (or can be subsequently) assigned to a different service for other transport protocols.

The most critical issue in picking a number is selecting the desired range, i.e., System vs. User port numbers. The distinction was intended to indicate a difference in privilege; originally, System port numbers required privileged ('root') access, while User port numbers did not. That distinction has since blurred because some current systems do not limit access control to System port numbers and because some System services have been replicated on User numbers (e.g., IRC). Even so, System port number assignments have continued at an average rate of 3-4 per year over the past 7 years (2007-2013), indicating that the desire to keep this distinction continues.

As a result, the difference between System and User port numbers needs to be treated with caution. Developers are advised to treat services as if they are always run without privilege.

Even when developers seek a System port number assignment, it may be very difficult to obtain. System port number assignment requires IETF Review or IESG Approval and justification that both User and

Dynamic port number ranges are insufficient [RFC6335]. Thus this document recommends both:

>> Developers SHOULD NOT apply for System port number assignments because the increased privilege they are intended to provide is not always enforced.

>> System implementers SHOULD enforce the need for privilege for processes to listen on System port numbers.

At some future date, it might be useful to deprecate the distinction between System and User port numbers altogether. Services typically require elevated ('root') privileges to bind to a System port number, but many such services go to great lengths to immediately drop those privileges just after connection or other association establishment to reduce the impact of an attack using their capabilities. Such services might be more securely operated on User port numbers than on System port numbers. Further, if System port numbers were no longer assigned, as of 2014 it would cost only 180 of the 1024 System values (17%), or 180 of the overall 49152 assigned (System and User) values (<0.04%).

7.4. Support for Security

Just as a service is a way to obtain information or processing from a host over a network, a service can also be the opening through which to compromise that host. Protecting a service involves security, which includes integrity protection, source authentication, privacy, or any combination of these capabilities. Security can be provided in a number of ways, and thus:

>> New services SHOULD support security capabilities, either directly or via a content protection such as TLS [RFC5246] or DTLS [RFC6347] or transport protection such as TCP-AO [RFC5925]. Insecure versions of new or existing secure services SHOULD be avoided because of the new vulnerability they create.

Secure versions of legacy services that are not already security-capable via in-band negotiations can be very useful. However, there is no IETF consensus on when separate ports should be used for secure and insecure variants of the same service [RFC2595] [RFC2817] [RFC6335]. The overall preference is for use of a single port, as noted in Section 6 of this document and Section 7.2 of [RFC6335], but the appropriate approach depends on the specific characteristics of the service. As a result:

>> When requesting both secure and insecure port assignments for the same service, justification is expected for the utility and safety of each port as an independent service (Section 6). Precedent (e.g., citing other protocols that use a separate insecure port) is inadequate justification by itself.

It's also important to recognize that port number assignment is not itself a guarantee that traffic using that number provides the corresponding service, or that a given service is always offered only on its assigned port number. Port numbers are ultimately meaningful only between endpoints and any service can be run on any port. Thus:

>> Security SHOULD NOT rely on assigned port number distinctions alone; every service, whether secure or not, is likely to be attacked.

Applications for a new service that requires both a secure and insecure port may be found, on expert review, to be unacceptable, and may not be approved for allocation. Similarly, an application for a new port to support an insecure variant of an existing secure protocol may be found unacceptable. In both cases, the resulting security of the service in practice will be a significant consideration in the decision as to whether to assign an insecure port.

7.5. Support for Future Versions

Requests for assigned port numbers are expected to support multiple versions on the same assigned port number [RFC6335]. Versions are typically indicated in-band, either at the beginning of a connection or other association, or in each protocol message.

>> Version support SHOULD be included in new services rather than relying on different port number assignments for different versions.

>> Version numbers SHOULD NOT be included in either the service name or service description, to avoid the need to make additional port number assignments for future variants of a service.

Again, the assigned port number space is far too limited to be used as an indicator of protocol version or message type. Although this has happened in the past (e.g., for NFS), it should be avoided in new requests.

7.6. Transport Protocols

IANA assigns port numbers specific to one or more transport protocols, typically UDP [RFC768] and TCP [RFC793], but also SCTP [RFC4960], DCCP [RFC4340], and any other standard transport protocol. Originally, IANA port number assignments were concurrent for both UDP and TCP, and other transports were not indicated. However, to conserve the assigned port number space and to reflect increasing use of other transports, assignments are now specific only to the transport being used.

In general, a service should request assignments for multiple transports using the same service name and description on the same port number only when they all reflect essentially the same service. Good examples of such use are DNS and NFS, where the difference between the UDP and TCP services are specific to supporting each transport. E.g., the UDP variant of a service might add sequence numbers and the TCP variant of the same service might add in-band message delimiters. This document does not describe the appropriate selection of a transport protocol for a service.

>> Service names and descriptions for multiple transport port number assignments SHOULD match only when they describe the same service, excepting only enhancements for each supported transport.

When the services differ, it may be acceptable or preferable to use the same port number, but the service names and descriptions should be different for each transport/service pair, reflecting the differences in the services. E.g., if TCP is used for the basic control protocol and UDP for an alarm protocol, then the services might be "name-ctl" and "name-alarm". A common example is when TCP is used for a service and UDP is used to determine whether that service is active (e.g., via a unicast, broadcast, or multicast test message) [RFC1122]. IANA has, for several years, used the suffix "-disc" in service names to distinguish discovery services, such as are used to identify endpoints capable of a given service:

>> Names of discovery services SHOULD use an identifiable suffix; the suggestion is "-disc".

Some services are used for discovery, either in conjunction with a TCP service or as a stand-alone capability. Such services will be more reliable when using multicast rather than broadcast (over IPv4) because IP routers do not forward "all nodes" broadcasts (all 1's, i.e., 255.255.255.255 for IPv4) and have not been required to support subnet-directed broadcasts since 1999 [RFC1812] [RFC2644].

This issue is relevant only for IPv4 because IPv6 does not support broadcast.

>> UDP over IPv4 multi-host services SHOULD use multicast rather than broadcast.

Designers should be very careful in creating services over transports that do not support congestion control or error recovery, notably UDP. There are several issues that should be considered in such cases, as summarized in Table 1 in [RFC5405]. In addition, the following recommendations apply to service design:

>> Services that use multipoint communication SHOULD be scalable, and SHOULD NOT rely solely on the efficiency of multicast transmission for scalability.

>> Services SHOULD NOT use UDP as a performance enhancement over TCP, e.g., to circumnavigate TCP's congestion control.

7.7. When to Request an Assignment

Assignments are typically requested when a user has enough information to reasonably answer the questions in the IANA application. IANA applications typically take up to a few weeks to process, with some complex cases taking up to a month. The process typically involves a few exchanges between the IANA Ports Expert Review team and the applicant.

An application needs to include a description of the service, as well as to address key questions designed to help IANA determine whether the assignment is justified. The application should be complete and not refer solely to the Internet Draft, RFC, a website, or any other external documentation.

Services that are independently developed can be requested at any time, but are typically best requested in the last stages of design and initial experimentation, before any deployment has occurred that cannot easily be updated.

>> Users MUST NOT deploy implementations that use assigned port numbers prior their assignment by IANA.

>> Users MUST NOT deploy implementations that default to using the experimental System port numbers (1021 and 1022 [RFC4727]) outside a controlled environment where they can be updated with a subsequent assigned port [RFC3692].

Deployments that use unassigned port numbers before assignment complicate IANA management of the port number space. Keep in mind that this recommendation protects existing assignees, users of current services, and applicants for new assignments; it helps ensure that a desired number and service name are available when assigned. The list of currently unassigned numbers is just that - **currently** unassigned. It does not reflect pending applications. Waiting for an official IANA assignment reduces the chance that an assignment request will conflict with another deployed service.

Applications made through Internet Draft / RFC publication (in any stream) typically use a placeholder ("PORTNUM") in the text, and implementations use an experimental port number until a final assignment has been made [RFC6335]. That assignment is initially indicated in the IANA Considerations section of the document, which is tracked by the RFC Editor. When a document has been approved for publication, that request is forwarded to IANA for handling. IANA will make the new assignment accordingly. At that time, IANA may also request that the applicant fill out the application form on their website, e.g., when the RFC does not directly address the information expected as per [RFC6335]. "Early" assignments can be made when justified, e.g., for early interoperability testing, according to existing process [RFC7120] [RFC6335].

>> Users writing specifications SHOULD use symbolic names for port numbers and service names until an IANA assignment has been completed. Implementations SHOULD use experimental port numbers during this time, but those numbers MUST NOT be cited in documentation except as interim.

7.8. Squatting

"Squatting" describes the use of a number from the assignable range in deployed software without IANA assignment for that use, regardless of whether the number has been assigned or remains available for assignment. It is hazardous because IANA cannot track such usage and thus cannot avoid making legitimate assignments that conflict with such unauthorized usage.

Such "squatted" port numbers remain unassigned, and IANA retains the right to assign them when requested by other applicants. Application and service designers are reminded that it is never appropriate to use port numbers that have not been directly assigned [RFC6335]. In particular, any unassigned code from the assigned ranges will be assigned by IANA, and any conflict will be easily resolved as the protocol designer's fault once that happens (because they would not be the assignee). This may reflect in the public's judgment on the

quality of their expertise and cooperation with the Internet community.

Regardless, there are numerous services that have squatted on such numbers that are in widespread use. Designers who are using such port numbers are encouraged to apply for an assignment. Note that even widespread de-facto use may not justify a later IANA assignment of that value, especially if either the value has already been assigned to a legitimate applicant or if the service would not qualify for an assignment of its own accord.

7.9. Other Considerations

As noted earlier, System port numbers should be used sparingly, and it is better to avoid them altogether. This avoids the potentially incorrect assumption that the service on such port numbers run in a privileged mode.

Assigned port numbers are not intended to be changed; this includes the corresponding service name. Once deployed, it can be very difficult to recall every implementation, so the assignment should be retained. However, in cases where the current assignee of a name or number has reasonable knowledge of the impact on such uses, and is willing to accept that impact, the name or number of an assignment can be changed [RFC6335]

Aliases, or multiple service names for the same assigned port number, are no longer considered appropriate [RFC6335].

8. Security Considerations

This document focuses on the issues arising when designing services that require new port assignments. Section 7.4 addresses the security and security-related issues of that interaction.

When designing a secure service, the use of TLS [RFC5246], DTLS [RFC6347], or TCP-AO [RFC5925] mechanisms that protect transport protocols or their contents is encouraged. It may not be possible to use IPsec [RFC4301] in similar ways because of the different relationship between IPsec and port numbers and because applications may not be aware of IPsec protections.

This document reminds application and service designers that port numbers do not protect against denial of service attack or guarantee that traffic should be trusted. Using assigned numbers for port filtering isn't a substitute for authentication, encryption, and integrity protection. The port number alone should not be used to

avoid denial of service attacks or to manage firewall traffic because the use of port numbers is not regulated or validated.

The use of assigned port numbers is the antithesis of privacy because they are intended to explicitly indicate the desired application or service. Strictly, port numbers are meaningful only at the endpoints, so any interpretation elsewhere in the network can be arbitrarily incorrect. However, those numbers can also expose information about available services on a given host. This information can be used by intermediate devices to monitor and intercept traffic as well as to potentially identify key endpoint software properties ("fingerprinting"), which can be used to direct other attacks.

9. IANA Considerations

The entirety of this document focuses on suggestions that help ensure the conservation of port numbers and provide useful hints for issuing informative requests thereof.

10. References

10.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2780] Bradner, S., and V. Paxson, "IANA Allocation Guidelines For Values In the Internet Protocol and Related Headers", BCP 37, RFC 2780, March 2000.
- [RFC3692] Narten, T., "Assigning Experimental and Testing Numbers Considered Useful", BCP 82, RFC 3962, Jan. 2004.
- [RFC4727] Fenner, B., "Experimental Values in IPv4, IPv6, ICMPv4, ICMPv6, UDP, and TCP Headers", RFC 4727, November 2006.
- [RFC5246] Dierks, T., and E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.2", RFC 5246, August 2008.
- [RFC5405] Eggert, L., and G. Fairhurst, "Unicast UDP Usage Guidelines for Application Designers", BCP 145, RFC 5405, Nov. 2008.
- [RFC5925] Touch, J., Mankin, A., and R. Bonica, "The TCP Authentication Option", RFC 5925, June 2010.

- [RFC6335] Cotton, M., L. Eggert, J. Touch, M. Westerlund, and S. Cheshire, "Internet Assigned Numbers Authority (IANA) Procedures for the Management of the Service Name and Transport Protocol Port Number Registry", BCP 165, RFC 6335, August 2011.
- [RFC6347] Rescorla, E., and N. Modadugu, "Datagram Transport Layer Security Version 1.2", RFC 6347, January 2012.

10.2. Informative References

- [IEN112] Postel, J., "Transmission Control Protocol", IEN 112, August 1979.
- [RFC33] Crocker, S., "New Host-Host Protocol", RFC 33 February 1970.
- [RFC37] Crocker, S., "Network Meeting Epilogue", RFC 37, March 1970.
- [RFC38] Wolfe, S., "Comments on Network Protocol from NWG/RFC #36", RFC 38, March 1970.
- [RFC48] Postel, J., and S. Crocker, "Possible protocol plateau", RFC 48, April 1970.
- [RFC61] Walden, D., "Note on Interprocess Communication in a Resource Sharing Computer Network", RFC 61, July 1970.
- [RFC76] Bouknight, J., J. Madden, and G. Grossman, "Connection by name: User oriented protocol", RFC 76, October 1970.
- [RFC333] Bressler, R., D. Murphy, and D. Walden. "Proposed experiment with a Message Switching Protocol", RFC 333, May 1972.
- [RFC739] Postel, J., "Assigned numbers", RFC 739, November 1977.
- [RFC758] Postel, J., "Assigned numbers", RFC 758, August 1979.
- [RFC768] Postel, J., "User Datagram Protocol", RFC 768, August 1980.
- [RFC793] Postel, J., "Transmission Control Protocol" RFC 793, September 1981
- [RFC820] Postel, J., "Assigned numbers", RFC 820, August 1982.

- [RFC900] Reynolds, J., and J. Postel, "Assigned numbers", RFC 900, June 1984.
- [RFC959] Postel, J., and J. Reynolds, "FILE TRANSFER PROTOCOL (FTP)", RFC 959, October 1985.
- [RFC1122] Braden, B. (Ed.), "Requirements for Internet Hosts -- Communication Layers", RFC 1122, October 1989.
- [RFC1340] Reynolds, J., and J. Postel, "Assigned numbers", RFC 1340, July 1992.
- [RFC1700] Reynolds, J., and J. Postel, "Assigned numbers", RFC 1700, October 1994.
- [RFC1812] Baker, F. (Ed.), "Requirements for IP Version 4 Routers", RFC 1812, June 1995.
- [RFC1833] Srinivasan, R., "Binding Protocols for ONC RPC Version 2", RFC 1833, August 1995.
- [RFC2595] Newman, C., "Using TLS with IMAP, POP3 and ACAP", RFC 2595, June 1999.
- [RFC2644] Senie, D., "Changing the Default for Directed Broadcasts in Routers", RFC 2644, August 1999.
- [RFC2817] Khare, R., and S. Lawrence, "Upgrading to TLS Within HTTP/1.1", RFC 2817, May 2000.
- [RFC3232] Reynolds, J. (Ed.), "Assigned Numbers: RFC 1700 is Replaced by an On-line Database", RFC 3232, January 2002.
- [RFC3261] Rosenberg, J., H. Schulzrinne, G. Camarillo, A. Johnston, J. Peterson, R. Sparks, M. Handley, and E. Schooler, "SIP: Session Initiation Protocol", RFC 3261, June 2002.
- [RFC4301] Kent, S., and K. Seo, "Security Architecture for the Internet Protocol", RFC 4301, December 2005.
- [RFC4340] Kohler, E., M. Handley, and S. Floyd, "Datagram Congestion Control Protocol (DCCP)", RFC 4340, March 2006.
- [RFC4960] Stewart, R. (Ed.), "Stream Control Transmission Protocol", RFC 4960, September 2007.

- [RFC5245] Rosenberg, J., "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols", RFC 5245, April 2010.
- [RFC5389] Rosenberg, J., R. Mahy, P. Matthews, and D. Wing, "Session Traversal Utilities for NAT", RFC 5389, October 2008.
- [RFC5766] Mahy, R., P. Matthews, and J. Rosenberg, "Traversal Using Relays around NAT (TURN): Relay Extensions to Session Traversal Utilities for NAT (STUN)", RFC 5766, April 2010.
- [RFC6066] Eastlake 3rd, D., "Transport Layer Security (TLS) Extensions: Extension Definitions", RFC 6066, January 2011.
- [RFC6762] Cheshire, S., and M. Krochmal, "Multicast DNS", RFC 6762, February 2013.
- [RFC6763] Cheshire, S., and M. Krochmal, "DNS-Based Service Discovery", RFC 6763, February 2013.
- [RFC7120] Cotton, M., "Early IANA Allocation of Standards Track Code Points", BCP 100, RFC 7120, January 2014.
- [RFC7230] Fielding, R., (Ed.), and J. Reshke, (Ed.), "Hypertext Transfer Protocol (HTTP/1.1): Message Syntax and Routing", RFC 7230, June 2014.

11. Acknowledgments

This work benefitted from the feedback from David Black, Lars Eggert, Gorry Fairhurst, and Eliot Lear, as well as discussions of the IETF TSVWG WG.

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

Authors' Addresses

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

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

Network WG
Internet-Draft
Expires: January 4, 2015
Intended Status: Standards Track
Updates: RFC 2872 (if accepted)

James Polk
Subha Dhesikan
Cisco Systems
July 4, 2014

Resource Reservation Protocol (RSVP) Application-ID
Profiles for Voice and Video Streams
draft-ietf-tsvwg-rsvp-app-id-vv-profiles-02

Abstract

RFC 2872 defines an Resource Reservation Protocol (RSVP) object for application identifiers. This document uses that App-ID and gives implementers specific guidelines for differing voice and video stream identifications to nodes along a reservation path, creating specific profiles for voice and video session identification.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on July 4, 2014.

Copyright Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1.	Introduction	2
2.	RSVP Application-ID Template	3
3.	The Voice and Video Application-ID Profiles	4
3.1	The Broadcast video Profile	4
3.2	The Real-time Interactive Profile	5
3.3	The Multimedia Conferencing Profile	5
3.4	The Multimedia Streaming Profile	6
3.5	The Conversational Profile	6
4.	Security considerations	7
5.	IANA considerations	7
5.1	Application Profiles	7
5.1.1	Broadcast Profiles IANA Registry	8
5.1.2	Realtime-Interactive Profiles IANA Registry	8
5.1.3	Multimedia-Conferencing Profiles IANA Registry	9
5.1.4	Multimedia-Streaming Profiles IANA Registry	10
5.1.5	Conversational Profiles IANA Registry	10
6.	Acknowledgments	12
7.	References	12
7.1.	Normative References	12
7.2.	Informative References	13
	Authors' Addresses	13
	Appendix	14

1. Introduction

RFC 2872 [RFC2872] describes the usage of policy elements for providing application information in Resource Reservation Protocol (RSVP) signaling [RFC2205]. The intention of providing this information is to enable application-based policy control. However, RFC 2872 does not enumerate any application profiles. The absence of explicit, uniform profiles leads to incompatible handling of these values and misapplied policies. An application profile used by a sender might not be understood by the intermediaries or receiver in a different domain. Therefore, there is a need to enumerate application profiles that are universally understood and applied for correct policy control.

Call control between endpoints has the ability to bind or associate many attributes to a reservation. One new attribute is currently being defined so as to establish the type of traffic contained in that reservation. This is accomplished via assigning a traffic label to the call (or session or flow) [ID-TRAF-CLASS].

This document takes the application traffic classes from [ID-TRAF-CLASS] and places those strings in the APP-ID object defined in RFC 2872. Thus, the intermediary devices (e.g., routers) processing the RSVP message can learn the identified profile within the Application-ID policy element for a particular reservation, and possibly be configured with the profile(s) to understand them

correctly, thus performing the correct admission control.

Another goal of this document is to the ability to signal an application profile which can then be translated into a DSCP value as per the choice of each domain. While the DCLASS object [RFC2996] allows the transfer of DSCP value in an RSVP message, that RFC does not allow the flexibility of having different domains choosing the DSCP value for the traffic classes that they maintain.

How these labels indicate the appropriate Differentiated Services Codepoint (DSCP) is out of scope for this document.

This document will break out each application type and propose how the values in application-id template should be populated for uniformity and interoperability.

1.1 Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC 2119].

2. RSVP Application ID Template

The template from RFC 2872 is as follows:

0	1	2	3
PE Length (8)	P-type = AUTH_APP		
Attribute Length	A-type = POLICY_LOCATOR	Sub-type = ASCII_DN	
Application name as ASCII string (e.g. SAP.EXE)			

In line with how this policy element is constructed in RFC 2872, the A-type will remain "POLICY_LOCATOR".

The P-type field is first created in [RFC2752]. This document uses the existing P-type "AUTH_APP" for application traffic class.

The first Sub-type will be mandatory for every profile within this document, and will be "ASCII_DN". No other Sub-types are defined by any profile within this document, but MAY be included by individual implementations - and MUST be ignored if not understood by receiving implementations along the reservation path.

RFC 2872 states the #1 sub-element from RFC 2872 as the "identifier that uniquely identifies the application vendor", which is optional to include. This document modifies this vendor limitation so that the identifier need only be unique - and not limited to an application vendor (identifier). For example, this specification now allows an RFC that defines an industry recognizable term or string to be a valid identifier. For example, a term or string taken from another IETF document, such as "conversational" or "avconf" from [ID-TRAF-CLASS]. This sub-element is still optional to include.

The following subsections will define the values within the above template into specific profiles for voice and video identification.

3. The Voice and Video Application-ID Profiles

This section contains the elements of the Application ID policy object which is used to signal the application classes defined in [ID-TRAF-CLASS].

3.1 The Broadcast Profiles

Broadcast profiles are for minimally buffered one-way streaming flows, such as video surveillance, or Internet based concerts or non-VOD TV broadcasts such as live sporting events.

This document creates Broadcast profiles for

- Broadcast IPTV for audio and video
- Broadcast Live-events for audio and video
- Broadcast Surveillance for audio and video

Here is an example profile for identifying Broadcast Video-Surveillance

```
AUTH_APP, POLICY_LOCATOR, ASCII_DN,  
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,  
APP=broadcast.video.surveillance, VER="
```

[Editor's Note: "rfcXXXX" will be replaced with the RFC number assigned to the [ID-TRAF-CLASS] reference. This 'note' should be removed during the RFC-Editor review process.]

Where the Globally Unique Identifier (GUID) indicates the documented reference that created this well-known string [ID-TRAF-CLASS], the APP is the profile name with no spaces, and the "VER=" is included, but has no value at this time.

3.2 The Realtime Interactive Profiles

Realtime Interactive profiles are for on-line gaming, and both remote and virtual avconf applications, in which the timing is particularly important towards the feedback to uses of these applications. This traffic type will generally not be UDP based, with minimal tolerance to RTT delays.

This document creates Realtime Interactive profiles for

- Realtime-Interactive Gaming
- Realtime-Interactive Remote-Desktop
- Realtime-Interactive Virtualized-Desktop

Here is the profile for identifying Realtime-Interactive Gaming

```
AUTH_APP, POLICY_LOCATOR, ASCII_DN,  
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,  
APP=realtime-interactive.gaming, VER="
```

Where the Globally Unique Identifier (GUID) indicates the documented reference that created this well-known string [ID-TRAF-CLASS], the APP is the profile name with no spaces, and the "VER=" is included, but has no value, but MAY if versioning becomes important.

3.3 The Multimedia Conferencing Profiles

There will be Multimedia Conferencing profiles for presentation data, application sharing and whiteboarding, where these applications will most often be associated with a larger Conversational (audio and/or audio/video) conference. Timing is important, but some minimal delays are acceptable, unlike the case for Realtime-Interactive traffic.

This document creates Multimedia-Conferencing profiles for

- Multimedia-Conferencing presentation-data
- Multimedia-Conferencing presentation-video
- Multimedia-Conferencing presentation-audio
- Multimedia-Conferencing application-sharing
- Multimedia-Conferencing whiteboarding

Here is the profile for identifying Multimedia-Conferencing Application-sharing

```
AUTH_APP, POLICY_LOCATOR, ASCII_DN,  
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,  
APP=multimedia-conferencing.application-sharing, VER="
```

Where the Globally Unique Identifier (GUID) indicates the RFC reference that created this well-known string [ID-TRAF-CLASS], the

APP is the profile name with no spaces, and the "VER=" is included, but has no value, but MAY if versioning becomes important.

3.4 The Multimedia Streaming Profiles

Multimedia Streaming profiles are for more significantly buffered one-way streaming flows than Broadcast profiles. These include...

This document creates Multimedia Streaming profiles for

- Multimedia-Streaming multiplex
- Multimedia-Streaming webcast

Here is the profile for identifying Multimedia Streaming webcast

```
AUTH_APP, POLICY_LOCATOR, ASCII_DN,  
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,  
APP=multimedia-streaming.webcast, VER="
```

Where the Globally Unique Identifier (GUID) indicates the documented reference that created this well-known string [ID-TRAF-CLASS], the APP is the profile name with no spaces, and the "VER=" is included, but has no value, but MAY if versioning becomes important.

3.5 The Conversational Profiles

Conversational category is for realtime bidirectional communications, such as voice or video, and is the most numerous due to the choices of application with or without adjectives. The number of profiles is then doubled because there needs to be one for unadmitted and one for admitted. The IANA section lists all that are currently proposed for registration at this time, therefore there will not be an exhaustive list provided in this section.

This document creates Conversational profiles for

- Conversational Audio
- Conversational Audio Admitted
- Conversational Video
- Conversational Video Admitted
- Conversational Audio Avconf
- Conversational Audio Avconf Admitted
- Conversational Video Avconf
- Conversational Video Avconf Admitted
- Conversational Audio Immersive
- Conversational Audio Immersive Admitted
- Conversational Video Immersive
- Conversational Video Immersive Admitted

Here is an example profile for identifying Conversational Audio:

```
AUTH_APP, POLICY_LOCATOR, ASCII_DN,  
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,  
APP=conversational.audio, VER="
```

Where the Globally Unique Identifier (GUID) indicates the documented reference that created this well-known string [ID-TRAF-CLASS], the APP is the profile name with no spaces, and the "VER=" is included, but has no value, but MAY if versioning becomes important.

4. Security considerations

The security considerations section within RFC 2872 sufficiently covers this document, with one possible exception - someone using the wrong template values (e.g., claiming a reservation is Multimedia Streaming when it is in fact Real-time Interactive). Given that each traffic flow is within separate reservations, and RSVP does not have the ability to police the type of traffic within any reservation, solving for this appears to be administratively handled at best. This is not meant to be a 'punt', but there really is nothing this template creates that is going to make things any harder for anyone (that we know of now).

5. IANA considerations

5.1 Application Profiles

This document requests IANA create a new registry for the application identification classes similar to the following table within the Resource Reservation Protocol (RSVP) Parameters registry:

```
Registry Name: RSVP APP-ID Profiles  
Reference: [this document]  
Registration procedures: Standards Track document [RFC5226]
```

```
[Editor's Note: "rfcXXXX" will be replaced with the RFC number  
assigned to the [ID-TRAF-CLASS] reference. This  
'note' should be removed during the RFC-Editor  
review process.]
```

5.1.1 Broadcast Profiles IANA Registry

Broadcast Audio IPTV Profile

```
P-type = AUTH_APP  
A-type = POLICY_LOCATOR  
Sub-type = ASCII_DN  
Conformant policy locator =  
    "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,  
    APP=broadcast.audio.iptv, VER="
```

Reference: [this document]

Broadcast Video IPTV Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=broadcast.video.iptv, VER="

Reference: [this document]

Broadcast Audio Live-events Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=broadcast.audio.live-events, VER="

Reference: [this document]

Broadcast Video Live-events Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=broadcast.video.live-events, VER="

Reference: [this document]

Broadcast Audio-Surveillance Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=broadcast.audio.surveillance, VER="

Reference: [this document]

Broadcast Video-Surveillance Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=broadcast.video.surveillance, VER="

Reference: [this document]

5.1.2 Realtime-Interactive Profiles IANA Registry

Realtime-Interactive Gaming Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP= realtime-interactive.gaming, VER="

Reference: [this document]

Real-time Interactive Remote-Desktop Profile

P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=realtime-interactive.remote-desktop, VER="

Reference: [this document]

Real-time Interactive Virtualized-Desktop Profile

P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=realtime-interactive.
remote-desktop.virtual, VER="

Reference: [this document]

Real-time Interactive Telemetry Profile

P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=realtime-interactive.telemetry, VER="

Reference: [this document]

5.1.3 Multimedia-Conferencing Profiles IANA Registry

Multimedia-Conferencing Presentation-Data Profile

P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP= multimedia-conferencing.presentation-data,
VER="

Reference: [this document]

Multimedia-Conferencing Presentation-Video Profile

P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,

APP= multimedia-conferencing.presentation-video,
VER="

Reference: [this document]

Multimedia-Conferencing Presentation-Audio Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP= multimedia-conferencing.presentation-audio,
VER="

Reference: [this document]

Multimedia-Conferencing Application-Sharing Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP= multimedia-conferencing.application-sharing,
VER="

Reference: [this document]

Multimedia-Conferencing Whiteboarding Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP= multimedia-conferencing.whiteboarding, VER="

Reference: [this document]

5.1.4 Multimedia-Streaming Profiles IANA Registry

Multimedia-Streaming Multiplex Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=multimedia-streaming.multiplex, VER="

Reference: [this document]

Multimedia-Streaming Webcast Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=multimedia-streaming.webcast, VER="

Reference: [this document]

5.1.5 Conversational Profiles IANA Registry

Conversational Audio Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=conversational.audio, VER="

Reference: [this document]

Conversational Audio Admitted Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=conversational.audio.aq:admitted, VER="

Reference: [this document]

Conversational Video Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=conversational.video, VER="

Reference: [this document]

Conversational Video Admitted Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=conversational.video.aq:admitted, VER="

Reference: [this document]

Conversational Audio Avconf Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=conversational.audio.avconf, VER="

Reference: [this document]

Conversational Audio Avconf Admitted Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
 "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
 APP=conversational.audio.avconf.aq:admitted,
 VER="

Reference: [this document]

Conversational Video Avconf Profile
P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
 "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
 APP=conversational.video.avconf, VER="

Reference: [this document]

Conversational Video Avconf Admitted Profile
P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
 "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
 APP=conversational.video.avconf.aq:admitted,
 VER="

Reference: [this document]

Conversational Audio Immersive Profile
P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
 "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
 APP=conversational.audio.immersive, VER="

Reference: [this document]

Conversational Audio Immersive Admitted Profile
P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
 "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
 APP=conversational.audio.immersive.aq:admitted,
 VER="

Reference: [this document]

Conversational Video Immersive Profile
P-type = AUTH_APP
A-type = POLICY_LOCATOR
Sub-type = ASCII_DN
Conformant policy locator =
 "GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,

APP=conversational.video.immersive, VER="

Reference: [this document]

Conversational Video Immersive Admitted Profile

P-type = AUTH_APP

A-type = POLICY_LOCATOR

Sub-type = ASCII_DN

Conformant policy locator =

"GUID=http://www.rfc-editor.org/rfc/rfcXXXX.txt,
APP=conversational.video.immersive.aq:admitted,
VER="

Reference: [this document]

6. Acknowledgments

To Francois Le Faucheur, Paul Jones, Ken Carlberg, Georgios Karagiannis and Glen Lavers for their helpful comments, document reviews and encouragement.

7. References

7.1. Normative References

- [RFC2119] S. Bradner, "Key words for use in RFCs to Indicate Requirement Levels", RFC 2119, March 1997
- [RFC2205] R. Braden, Ed., L. Zhang, S. Berson, S. Herzog, S. Jamin, "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification", RFC 2205, September 1997
- [RFC2474] K. Nichols, S. Blake, F. Baker, D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers ", RFC 2474, December 1998
- [RFC2750] S. Herzog, "RSVP Extensions for Policy Control", RFC 2750, January 2000
- [RFC2872] Y. Bernet, R. Pabbati, "Application and Sub Application Identity Policy Element for Use with RSVP", RFC 2872, June 2000
- [RFC2996] Y. Bernet, "Format of the RSVP DCLASS Object", RFC 2996, November 2000
- [RFC3182] S. Yadav, R. Yavatkar, R. Pabbati, P. Ford, T. Moore, S. Herzog, R. Hess, "Identity Representation for RSVP", RFC 3182, October 2001
- [RFC5226] T. Narten, H. Alvestrand, "Guidelines for Writing an IANA

Considerations Section in RFCs", RFC 5226, May 2008

[ID-TRAF-CLASS] J. Polk, S. Dhesikan, P. Jones, "The Session Description Protocol (SDP) 'trafficclass' Attribute", work in progress, Feb 2013

7.2. Informative References

[RFC4594] J. Babiarez, K. Chan, F Baker, "Configuration Guidelines for Diffserv Service Classes", RFC 4594, August 2006

Authors' Addresses

James Polk
3913 Treemont Circle
Colleyville, Texas, USA
+1.817.271.3552

mailto: jmpolk@cisco.com

Subha Dhesikan
170 W Tasman St
San Jose, CA, USA
+1.408-902-3351

mailto: sdhesika@cisco.com

Appendix - Changes to ID

[Editor's Note: this appendix should be removed in the RFC-Editor's process.]

A.1 - Changes from WG version -00 to WG version -01

The following changes were made in this version:

- corrected nits
- globally replaced GUID link from the MMUSIC Trafficclass ID to the future RFC of that document.
- added profiles for presentation-video and presentation-audio

A.2 - Changes from Individual -04 to WG version -00

The following changes were made in this version:

- changed P-Type from APP_TC back to AUTH_APP, which is already defined.
- fixed nits and inconsistencies

A.3 - Changes from Individual -03 to -04

The following changes were made in this version:

- clarified security considerations section to mean RSVP cannot police the type of traffic within a reservation to know if a traffic flow should be using a different profile, as defined in this document.
- changed existing informative language regarding "... other Sub-types ..." from 'can' to normative 'MAY'.
- editorial changes to clear up minor mistakes

A.4 - Changes from Individual -02 to -03

The following changes were made in this version:

- Added [ID-TRAF-CLASS] as a reference
- Changed to a new format of the profile string.
- Added many new profiles based on the new format into each parent category of Section 3.
- changed the GUID to refer to draft-ietf-mmusic-traffic-class-for-sdp-03.txt
- changed 'desktop' adjective to 'avconf' to keep in alignment with [ID-TRAF-CLASS]
- Have a complete IANA Registry proposal for each application-ID discussed in this draft.
- General text clean-up of the draft.

Internet Engineering Task Force
Internet-Draft
Intended status: Experimental
Expires: April 6, 2015

Georgios Karagiannis
Huawei Technologies
Anurag Bhargava
Cisco Systems, Inc.
October 6, 2014

Generic Aggregation of Resource ReSerVation Protocol (RSVP)
for IPv4 And IPv6 Reservations over PCN domains
draft-ietf-tsvwg-rsvp-pcn-11

Abstract

This document specifies extensions to Generic Aggregated RSVP RFC 4860 for support of the PCN Controlled Load (CL) and Single Marking (SM) edge behaviors over a Diffserv cloud using Pre-Congestion Notification.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on April 6, 2015.

Copyright Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC2119].

Table of Contents

1. Introduction	4
1.1. Objective	4
1.2. Overview and Motivation	5
1.3. Terminology	7
1.4. Organization of This Document	11
2. Overview of RSVP extensions and Operations	11
2.1. Overview of RSVP Aggregation Procedures in PCN domains	11
2.2. PCN Marking and encoding and transport of pre-congestion Information	13
2.3. Traffic Classification Within The Aggregation Region	13
2.4. Deaggregator (PCN-egress-node) Determination	13
2.5. Mapping E2E Reservations Onto Aggregate Reservations	13
2.6. Size of Aggregate Reservations	14
2.7. E2E Path ADSPEC update	14
2.8. Intra-domain Routes	14
2.9. Inter-domain Routes	15
2.10. Reservations for Multicast Sessions	15
2.11. Multi-level Aggregation	15
2.12. Reliability Issues	15
3. Elements of Procedure	15
3.1. Receipt of E2E Path Message by PCN-ingress-node (aggregating router)	15
3.2. Handling Of E2E Path Message by Interior Routers	16
3.3. Receipt of E2E Path Message by PCN-egress-node (deaggregating router)	16
3.4. Initiation of new Aggregate Path Message By PCN-ingress-node (Aggregating Router)	16
3.5. Handling Of new Aggregate Path Message by Interior Routers	16
3.6. Handling Of Aggregate Path Message by Deaggregating Router	16
3.7. Handling of E2E Resv Message by Deaggregating Router	17
3.8. Handling Of E2E Resv Message by Interior Routers	17

3.9. Initiation of New Aggregate Resv Message By Deaggregating Router	17
3.10. Handling of Aggregate Resv Message by Interior Routers	18
3.11. Handling of E2E Resv Message by Aggregating Router	18
3.12. Handling of Aggregated Resv Message by Aggregating Router . .	18
3.13. Removal of E2E Reservation	19
3.14. Removal of Aggregate Reservation	19
3.15. Handling of Data On Reserved E2E Flow by Aggregating Router .	19
3.16. Procedures for Multicast Sessions	19
3.17. Misconfiguration of PCN node	19
3.18. PCN based Flow Termination	19
4. Protocol Elements	20
4.1 PCN object	20
5. Security Considerations	23
6. IANA Considerations	24
7. Acknowledgments	24
8. Normative References	24
9. Informative References	25
10. Appendix A: Example Signaling Flow	26
11. Authors' Address	29

1. Introduction

1.1 Objective

Pre-Congestion Notification (PCN) can support the quality of service (QoS) of inelastic flows within a Diffserv domain in a simple, scalable, and robust fashion. Two mechanisms are used: admission control and flow termination. Admission control is used to decide whether to admit or block a new flow request, while flow termination is used in abnormal circumstances to decide whether to terminate some of the existing flows. To support these two features, the overall rate of PCN-traffic is metered on every link in the domain, and PCN-packets are appropriately marked when certain configured rates are exceeded. These configured rates are below the rate of the link, thus providing notification to boundary nodes about overloads before any congestion occurs (hence "pre-congestion" notification). The PCN-egress-nodes measure the rates of differently marked PCN traffic in periodic intervals and report these rates to the Decision Points for admission control and flow termination; the Decision Points use these rates to make decisions. The Decision Points may be collocated with the PCN-ingress-nodes, or their function may be implemented in a another node. For more details see [RFC5559], [RFC6661], and [RFC6662].

The main objective of this document is to specify the signaling protocol that can be used within a Pre-Congestion Notification (PCN) domain to carry reports from a PCN-ingress-node to a PCN Decision point, considering that the PCN Decision Point and PCN-egress-node are collocated.

If the PCN Decision Point is not collocated with the PCN-egress-node then additional signaling procedures are required that are out of the scope of this document. Moreover, as mentioned above this architecture conforms with PBAC (Policy-Based Admission Control), when the Decision Point is located in a another node then the PCN-ingress-node [RFC2753].

Several signaling protocols can be used to carry information between PCN-boundary-nodes (PCN-ingress-node and PCN-egress-node). However, since (1) both PCN-egress-node and PCN-ingress-nodes are located on the data path and (2) the admission control procedure needs to be done at PCN-egress-node, a signaling protocol that follows the same path as the data path, like RSVP (Resource Reservation Protocol), is more suited for this purpose. In particular, this document specifies extensions to Generic Aggregated RSVP [RFC4860] for support of the PCN Controlled Load (CL) and Single Marking (SM) edge behaviors over a Diffserv cloud using Pre-Congestion Notification.

This draft is intended to be published as Experimental in order to:

- o) validate industry interest by allowing implementation and deployment
- o) gather operational experience, in particular around dynamic interactions of RSVP signaling and PCN notification and

corresponding levels of performance.

Support for the techniques specified in this document involves RSVP functionality in boundary nodes of a PCN domain whose interior nodes forward RSVP traffic without performing RSVP functionality.

1.2 Overview and Motivation

Two main Quality of Service (QoS) architectures have been specified by the IETF. These are the Integrated Services (Intserv) [RFC1633] architecture and the Differentiated Services (DiffServ) architecture ([RFC2475]).

Intserv provides methods for the delivery of end-to-end Quality of Service (QoS) to applications over heterogeneous networks. One of the QoS signaling protocols used by the Intserv architecture is the Resource reSerVation Protocol (RSVP) [RFC2205], which can be used by applications to request per-flow resources from the network. These RSVP requests can be admitted or rejected by the network. Applications can express their quantifiable resource requirements using Intserv parameters as defined in [RFC2211] and [RFC2212]. The Controlled Load (CL) service [RFC2211] is a quality of service (QoS) closely approximating the QoS that the same flow would receive from a lightly loaded network element. The CL service is useful for inelastic flows such as those used for real-time media.

The DiffServ architecture can support the differentiated treatment of packets in very large scale environments. While Intserv and RSVP classify packets per-flow, Diffserv networks classify packets into one of a small number of aggregated flows or "classes", based on the Diffserv codepoint (DSCP) in the packet IP header. At each Diffserv router, packets are subjected to a "per-hop behavior" (PHB), which is invoked by the DSCP. The primary benefit of Diffserv is its scalability, since the need for per-flow state and per-flow processing, is eliminated.

However, DiffServ does not include any mechanism for communication between applications and the network. Several solutions have been specified to solve this issue. One of these solutions is Intserv over Diffserv [RFC2998] including resource-based admission control (RBAC), PBAC, assistance in traffic identification/classification, and traffic conditioning. Intserv over Diffserv can operate over a statically provisioned or a RSVP aware Diffserv region. When it is RSVP aware, several mechanisms may be used to support dynamic provisioning and topology-aware admission control, including aggregate RSVP reservations, per-flow RSVP, or a bandwidth broker. [RFC3175] specifies aggregation of Resource ReSerVation Protocol (RSVP) end-to-end reservations over aggregate RSVP reservations. In [RFC3175] the RSVP generic aggregated reservation is characterized by a RSVP SESSION object using the 3-tuple <source IP address, destination IP address, Diffserv Code Point>.

Several scenarios require the use of multiple generic aggregate reservations that are established for a given PHB from a given source

IP address to a given destination IP address, see [SIG-NESTED], [RFC4860]. For example, multiple generic aggregate reservations can be applied in the situation that multiple E2E reservations using different preemption priorities need to be aggregated through a PCN-domain using the same PHB. By using multiple aggregate reservations for the same PHB, it allows enforcement of the different preemption priorities within the aggregation region. This allows more efficient management of the Diffserv resources, and in periods of resource shortage, this allows sustainment of a larger number of E2E reservations with higher preemption priorities. In particular, [SIG-NESTED] discusses in detail how end-to-end RSVP reservations can be established in a nested VPN environment through RSVP aggregation.

[RFC4860] provides generic aggregate reservations by extending [RFC3175] to support multiple aggregate reservations for the same source IP address, destination IP address, and PHB (or set of PHBs). In particular, multiple such generic aggregate reservations can be established for a given PHB from a given source IP address to a given destination IP address. This is achieved by adding the concept of a Virtual Destination Port and of an Extended Virtual Destination Port in the RSVP SESSION object. In addition to this, the RSVP SESSION object for generic aggregate reservations uses the PHB Identification Code (PHB-ID) defined in [RFC3140], instead of using the Diffserv Code Point (DSCP) used in [RFC3175]. The PHB-ID is used to identify the PHB, or set of PHBs, from which the Diffserv resources are to be reserved.

The RSVP like signaling protocol required to carry (1) requests from a PCN-egress-node to a PCN-ingress-node and (2) reports from a PCN-ingress-node to a PCN-egress-node needs to follow the PCN signaling requirements defined in [RFC6663]. In addition to that the signaling protocol functionality supported by the PCN-ingress-nodes and PCN-egress-nodes needs to maintain logical aggregate constructs (i.e. ingress-egress-aggregate state) and be able to map E2E reservations to these aggregate constructs. Moreover, no actual reservation state is needed to be maintained inside the PCN domain, i.e., the PCN-interior-nodes are not maintaining any reservation state.

This can be accomplished by two possible approaches:

Approach (1):

- o) adapting the RFC 4860 aggregation procedures to fit the PCN requirements with as little change as possible over the RFC 4860 functionality
- o) hence performing aggregate RSVP signaling (even if it is to be ignored by PCN interior nodes)
- o) using this aggregate RSVP signaling procedures to carry PCN information between the PCN-boundary-nodes (PCN-ingress-node and PCN-egress-node).

Approach (2):

- o) adapting the RFC 4860 aggregation procedures to fit the PCN requirements with more significant changes over RFC4860 (i.e. the aspect of the procedures that have to do with maintaining aggregate states and to do with mapping the E2E reservations to aggregate constructs are kept, but the procedures that have to do with the aggregate RSVP signaling and aggregate reservation establishment/maintenance are dropped).
- o) hence not performing aggregate RSVP signaling
- o) piggy-backing of the PCN information inside the E2E RSVP signaling.

Both approaches are probably viable, however, since the RFC 4860 operations have been thoroughly studied and implemented, it can be considered that the RFC 4860 solution can better deal with the more challenging situations (rerouting in the PCN domain, failure of an PCN-ingress-node, failure of an PCN-egress-node, rerouting towards a different edge, etc.). This is the reason for choosing Approach (1) for the specification of the signaling protocol used to carry PCN information between the PCN-boundary-nodes (PCN-ingress-node and PCN-egress-node).

In particular, this document specifies extensions to Generic Aggregated RSVP [RFC4860] for support of the PCN Controlled Load (CL) and Single Marking (SM) edge behaviors over a Diffserv cloud using Pre-Congestion Notification.

This document follows the PCN signaling requirements defined in [RFC6663] and specifies extensions to Generic Aggregated RSVP [RFC4860] for support of PCN edge behaviors as specified in [RFC6661] and [RFC6662]. Moreover, this document specifies how RSVP aggregation can be used to setup and maintain: (1) Ingress Egress Aggregate (IEA) states at Ingress and Egress nodes and (2) generic aggregation of RSVP end-to-end RSVP reservations over PCN (Congestion and Pre-Congestion Notification) domains.

To comply with this specification, PCN-nodes MUST be able to support the functionality specified in [RFC5670], [RFC5559], [RFC6660], [RFC6661], [RFC6662]. Furthermore, the PCN-boundary-nodes MUST support the RSVP generic aggregated reservation procedures specified in [RFC4860] which are augmented with procedures specified in this document.

1.3. Terminology

This document uses terms defined in [RFC4860], [RFC3175], [RFC5559], [RFC5670], [RFC6661], [RFC6662].

For readability, a number of definitions from [RFC3175] as well as definitions for terms used in [RFC5559], [RFC6661], and [RFC6662] are provided here, where some of them are augmented with new meanings:

Aggregator	This is the process in (or associated with) the router at the ingress edge of the aggregation region (with respect to the end-to-end RSVP reservation) and behaving in accordance with [RFC4860]. In this document, it is also the PCN-ingress-node. It is important to notice that in the context of this document the Aggregator must be able to determine the Deaggregator using the procedures specified in Section 4 of [RFC4860] and in Section 1.4.2 of [RFC3175].
Congestion level estimate (CLE):	<p>The ratio of PCN-marked to total PCN-traffic (measured in octets) received for a given ingress-egress-aggregate during a given measurement period. The CLE is used to derive the PCN-admission-state and is also used by the report suppression procedure if report suppression is activated.</p>
Deaggregator	This is the process in (or associated with) the router at the egress edge of the aggregation region (with respect to the end-to-end RSVP reservation) and behaving in accordance with [RFC4860]. In this document, it is also the PCN-egress-node and Decision Point.
E2E	end to end
E2E Reservation	<p>This is an RSVP reservation such that:</p> <ul style="list-style-type: none">(i) corresponding RSVP Path messages are initiated upstream of the Aggregator and terminated downstream of the Deaggregator, and(ii) corresponding RSVP Resv messages are initiated downstream of the Deaggregator and terminated upstream of the Aggregator, and(iii) this RSVP reservation is aggregated over an Ingress Egress Aggregate (IEA) between the Aggregator and Deaggregator. <p>An E2E RSVP reservation may be a per-flow reservation, which in this document is only maintained at the PCN-ingress-node and PCN-egress-node. Alternatively, the E2E reservation may itself be an aggregate reservation of various types (e.g., Aggregate IP reservation, Aggregate IPsec reservation, see [RFC4860]). As per regular RSVP operations, E2E RSVP reservations are unidirectional.</p>
E2E microflow	a microflow where its associated packets are being forwarded on an E2E path.

Extended vDstPort (Extended Virtual Destination Port)

An identifier used in the SESSION that remains constant over the life of the generic aggregate reservation. The length of this identifier is 32-bits when IPv4 addresses are used and 128 bits when IPv6 addresses are used.

A sender(or Aggregator) that wishes to narrow the scope of a SESSION to the sender-receiver pair (or Aggregator-Deaggregator pair) should place its IPv4 or IPv6 address here as a network unique identifier. A sender (or Aggregator) that wishes to use a common session with other senders (or Aggregators) in order to use a shared reservation across senders (or Aggregators) must set this field to all zeros. In this document, the Extended vDstPort should contain the IPv4 or IPv6 address of the Aggregator.

ETM-rate

The rate of excess-traffic-marked PCN-traffic received at a PCN-egress-node for a given ingress-egress-aggregate in octets per second.

Ingress-egress-aggregate (IEA):

The collection of PCN-packets from all PCN-flows that travel in one direction between a specific pair of PCN-boundary-nodes. In this document one RSVP generic aggregated reservation is mapped to only one ingress-egress-aggregate, while one ingress-egress-aggregate is mapped to either one or to more than one RSVP generic aggregated reservations. PCN-flows and their PCN-traffic that are mapped into a specific RSVP generic aggregated reservation can also easily be mapped into their corresponding ingress-egress-aggregate.

Microflow:
(from [RFC2474])

a single instance of an application-to-application flow of packets which is identified by source address, destination address, protocol id, and source port, destination port (where applicable).

PCN-domain:

a PCN-capable domain; a contiguous set of PCN-enabled nodes that perform Diffserv scheduling [RFC2474]; the complete set of PCN-nodes that in principle can, through PCN-marking packets, influence decisions about flow admission and termination within the domain; includes the PCN-egress-nodes, which measure these PCN-marks, and the PCN-ingress-nodes.

PCN-boundary-node: a PCN-node that connects one PCN-domain to a node either in another PCN-domain or in a non-PCN-domain.

- PCN-interior-node: a node in a PCN-domain that is not a PCN-boundary-node.
- PCN-node: a PCN-boundary-node or a PCN-interior-node.
- PCN-egress-node: a PCN-boundary-node in its role in handling traffic as it leaves a PCN-domain. In this document the PCN-egress-node operates also as a Decision Point and Deaggregator.
- PCN-ingress-node: a PCN-boundary-node in its role in handling traffic as it enters a PCN-domain. In this document the PCN-ingress-node operates also as a Aggregator.
- PCN-traffic,
PCN-packets,
PCN-BA: a PCN-domain carries traffic of different Diffserv behavior aggregates (BAs) [RFC2474]. The PCN-BA uses the PCN mechanisms to carry PCN-traffic, and the corresponding packets are PCN-packets. The same network will carry traffic of other Diffserv BAs. The PCN-BA is distinguished by a combination of the Diffserv codepoint (DSCP) and ECN fields.
- PCN-flow: the unit of PCN-traffic that the PCN-boundary-node admits (or terminates); the unit could be a single E2E microflow (as defined in [RFC2474]) or some identifiable collection of microflows.
- PCN-admission-state: The state ("admit" or "block") derived by the Decision Point for a given ingress-egress-aggregate based on statistics about PCN-packet marking. The Decision Point decides to admit or block new flows offered to the aggregate based on the current value of the PCN-admission-state.
- PCN-sent-rate The rate of PCN-traffic received at a PCN-ingress-node and destined for a given ingress-egress-aggregate in octets per second.
- PHB-ID (Per Hop Behavior Identification Code)
A 16-bit field containing the Per Hop Behavior Identification Code of the PHB, or of the set of PHBs, from which Diffserv resources are to be reserved. This field must be encoded as specified in Section 2 of [RFC3140].
- RSVP generic aggregated reservation: an RSVP reservation that is identified by using the RSVP SESSION object for generic RSVP aggregated reservation. This RSVP

SESSION object is based on the RSVP SESSION object specified in [RFC4860] augmented with the following information:

- o) the IPv4 DestAddress, IPv6 DestAddress should be set to the IPv4 or IPv6 destination addresses, respectively, of the Deaggregator (PCN-egress-node)
- o) PHB-ID (Per Hop Behavior Identification Code) should be set equal to PCN-compatible Diffserv codepoint(s).
- o) Extended vDstPort should be set to the IPv4 or IPv6 destination addresses, of the Aggregator (PCN-ingress-node)

VDstPort (Virtual Destination Port)

A 16-bit identifier used in the SESSION that remains constant over the life of the generic aggregate reservation.

1.4. Organization of This Document

This document is organized as follows. Section 2 gives an overview of RSVP extensions and operations. The elements of the used procedures are specified in Section 3. Section 4 describes the protocol elements. The security considerations are given in section 5 and the IANA considerations are provided in Section 6.

2. Overview of RSVP extensions and Operations

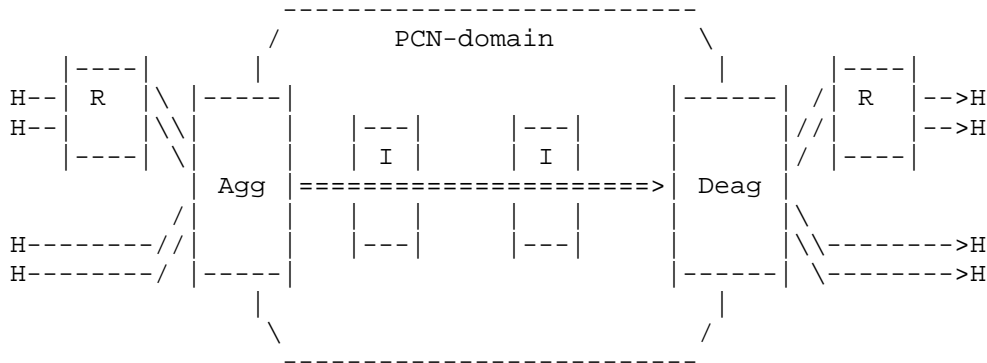
2.1 Overview of RSVP Aggregation Procedures in PCN domains

The PCN-boundary-nodes, see Figure 1, can support RSVP SESSIONS for generic aggregated reservations [RFC4860], which are depending on ingress-egress-aggregates. In particular, one RSVP generic aggregated reservation matches to only one ingress-egress-aggregate.

However, one ingress-egress-aggregate matches to either one, or more than one, RSVP generic aggregated reservations. In addition, to comply with this specification, the PCN-boundary nodes need to distinguish and process (1) RSVP SESSIONS for generic aggregated sessions and their messages according to [RFC4860], (2) E2E RSVP sessions and messages according to [RFC2205].

This document locates all RSVP processing for a PCN domain at PCN-Boundary nodes. PCN-interior-nodes do not perform any RSVP functionality or maintain RSVP-related state information. Rather, PCN-interior nodes forward all RSVP messages (for both generic aggregated reservations[RFC4860] and end to end reservations [RFC2205]) as if they were ordinary network traffic.

Moreover, each Aggregator and Deaggregator (i.e., PCN-boundary-nodes) need to support policies to initiate and maintain for each pair of PCN-boundary-nodes of the same PCN-domain one ingress-egress-aggregate.



H = Host requesting end-to-end RSVP reservations
 R = RSVP router
 Agg = Aggregator (PCN-ingress-node)
 Deag = Deaggregator (PCN-egress-node)
 I = Interior Router (PCN-interior-node)
 --> = E2E RSVP reservation
 ==> = Aggregate RSVP reservation

Figure 1 : Aggregation of E2E Reservations
 over Generic Aggregate RSVP Reservations
 in PCN domains, based on [RFC4860]

Both the Aggregator and Deaggregator can maintain one or more RSVP generic aggregated Reservations, but the Deaggregator is the entity that initiates these RSVP generic aggregated reservations. Note that one RSVP generic aggregated reservation matches to only one ingress-egress-aggregate, while one ingress-egress-aggregate matches to either one or to more than one RSVP generic aggregated reservations. This can be accomplished by using for the different RSVP generic aggregated reservations the same combinations of ingress and egress identifiers, but with a different PHB-ID value (see [RFC4860]). The procedures for aggregation of E2E reservations over generic aggregate RSVP reservations are the same as the procedures specified in Section 4 of [RFC4860], augmented with the ones specified in Section 2.5.

One significant difference between this document and [RFC4860] is the fact that in this document the admission control of E2E RSVP reservations over the PCN core is performed according to the PCN procedures, while in [RFC4860] this is achieved via first admitting aggregate RSVP reservations over the aggregation region and then admitting the E2E reservations over the aggregate RSVP reservations. Therefore, in this document, the RSVP generic aggregate RSVP reservations are not subject to admission control in the PCN-core, and the E2E RSVP reservations are not subject to admission control

over the aggregate reservations. In turn, this means that several procedures of [RFC4860] are significantly simplified in this document:

- o) unlike [RFC4860], the generic aggregate RSVP reservations need not be admitted in the PCN core.
- o) unlike [RFC4860], the RSVP aggregated traffic does not need to be tunneled between Aggregator and Deaggregator, see Section 2.3.
- o) unlike [RFC4860], the Deaggregator need not perform admission control of E2E reservations over the aggregate RSVP reservations.
- o) unlike [RFC4860], there is no need for dynamic adjustment of the RSVP generic aggregated reservation size, see Section 2.6.

2.2 PCN Marking and encoding and transport of pre-congestion information

The method of PCN marking within the PCN domain is specified in [RFC5670]. In addition, the method of encoding and transport of pre-congestion information is specified in [RFC6660]. The PHB-ID (Per Hop Behavior Identification Code) used SHOULD be set equal to PCN-compatible Diffserv codepoint(s).

2.3. Traffic Classification Within The Aggregation Region

The PCN-ingress marks a PCN-BA using PCN-marking (i.e., combination of the DSCP and ECN fields), which interior nodes use to classify PCN-traffic. The PCN-traffic (e.g., E2E microflows) belonging to a RSVP generic aggregated reservation can be classified only at the PCN-boundary-nodes (i.e., Aggregator and Deaggregator) by using the RSVP SESSION object for RSVP generic aggregated reservations, see Section 2.1 of [RFC4860]. Note that the DSCP value included in the SESSION object, SHOULD be set equal to a PCN-compatible Diffserv codepoint. Since no admission control procedures over the RSVP generic aggregated reservations in the PCN-core are required, unlike [RFC4860], the RSVP aggregated traffic need not to be tunneled between Aggregator and Deaggregator. In this document one RSVP generic aggregated reservation is mapped to only one ingress-egress-aggregate, while one ingress-egress-aggregate is mapped to either one or to more than one RSVP generic aggregated reservations. PCN-flows and their PCN-traffic that are mapped into a specific RSVP generic aggregated reservation can also easily be classified into their corresponding ingress-egress-aggregate. The method of traffic conditioning of PCN-traffic and non-PCN traffic and PHB configuration is described in [RFC6661] and [RFC6662].

2.4. Deaggregator Determination

The present document assumes the same dynamic Deaggregator determination method as used in [RFC4860].

2.5. Mapping E2E Reservations Onto Aggregate Reservations

To comply with this specification for the mapping of E2E reservations

onto aggregate reservations, the same methods MUST be used as the ones described in Section 4 of [RFC4860], augmented by the following rules:

- o) An Aggregator (also PCN-ingress-node in this document) or Deaggregator (also PCN-egress-node and Decision Point in this document) MUST use one or more policies to determine whether a RSVP generic aggregated reservation can be mapped into an ingress-Egress-aggregate. This can be accomplished by using for the different RSVP generic aggregated reservations the same combinations of ingress and egress identifiers, but with a different PHB-ID value (see [RFC4860]) corresponding to the PCN specifications. In particular, the RSVP SESSION object specified in [RFC4860] augmented with the following information:
 - o) the IPv4 DestAddress, IPv6 DestAddress MUST be set to the IPv4 or IPv6 destination addresses, respectively, of the Deaggregator (PCN-egress-node), see [RFC4860]. Note that the PCN-domain is considered as being only one RSVP hop (for Generic aggregated RSVP or E2E RSVP). This means that the next RSVP hop for the Aggregator in the downstream direction is the Deaggregator and the next RSVP hop for the Deaggregator in the upstream direction is the Aggregator.
 - o) PHB-ID (Per Hop Behavior Identification Code) SHOULD be set equal to PCN-compatible Diffserv codepoint(s).
 - o) Extended vDstPort SHOULD be set to the IPv4 or IPv6 destination addresses, of the Aggregator (PCN-ingress-node), see [RFC4860].

2.6. Size of Aggregate Reservations

Since:(i) no admission control of E2 reservations over the RSVP aggregated reservations is required, and (ii) no admission control of the RSVP aggregated reservation over the PCN core is required, the size of the generic aggregate reservation is irrelevant and can be set to any arbitrary value by the Deaggregator. The Deaggregator SHOULD set the value of a generic aggregate reservation to a null bandwidth. We also observe that there is no need for dynamic adjustment of the RSVP aggregated reservation size.

2.7. E2E Path ADSPEC update

To comply with this specification, for the update of the E2E Path ADSPEC, the same methods can be used as the ones described in [RFC4860].

2.8. Intra-domain Routes

The PCN-interior-nodes are neither maintaining E2E RSVP nor RSVP generic aggregation states and reservations. Therefore, intra-domain route changes will not affect intra-domain reservations since such reservations are not maintained by the PCN-interior-nodes.

Furthermore, it is considered that by configuration, the PCN-interior-nodes are not able to distinguish neither RSVP generic aggregated sessions and their associated messages [RFC4860], nor E2E RSVP sessions and their associated messages [RFC2205].

2.9. Inter-domain Routes

The PCN-charter scope precludes inter-domain considerations. However, for solving inter-domain routes changes associated with the operation of the RSVP messages, the same methods SHOULD be used as the ones described in [RFC4860] and in Section 1.4.7 of [RFC3175].

2.10. Reservations for Multicast Sessions

PCN does not consider reservations for multicast sessions.

2.11. Multi-level Aggregation

PCN does not consider multi-level aggregations within the PCN domain. Therefore, the PCN-interior-nodes are not supporting multi-level aggregation procedures. However, the Aggregator and Deaggregator SHOULD support the multi-level aggregation procedures specified in [RFC4860] and in Section 1.4.9 of [RFC3175].

2.12. Reliability Issues

To comply with this specification, for solving possible reliability issues, the same methods MUST be used as the ones described in Section 4 of [RFC4860].

3. Elements of Procedure

This section describes the procedures used to implement the aggregated RSVP procedure over PCN. It is considered that the procedures for aggregation of E2E reservations over generic aggregate RSVP reservations are same as the procedures specified in Section 4 of [RFC4860] except where a departure from these procedures is explicitly described in the present section. Please refer to [RFC4860] for all the below error cases:

- o) Incomplete message
- o) Unexpected objects

3.1. Receipt of E2E Path Message by Aggregating router

When the E2E Path message arrives at the exterior interface of the Aggregator, (also PCN-ingress-node in this document), then standard RSVP generic aggregation [RFC4860] procedures are used.

3.2. Handling Of E2E Path Message by Interior Routers

The E2E Path messages traverse zero or more PCN-interior-nodes. The PCN-interior-nodes receive the E2E Path message on an interior interface and forward it on another interior interface. It is considered that, by configuration, the PCN-interior-nodes ignore the E2E RSVP signaling messages [RFC2205]. Therefore, the E2E Path messages are simply forwarded as normal IP datagrams.

3.3. Receipt of E2E Path Message by Deaggregating router

When receiving the E2E Path message the Deaggregator (also PCN-egress-node and Decision Point in this document) performs the regular [RFC4860] procedures, augmented with the following rules:

- o) The Deaggregator MUST NOT perform the RSVP-TTL vs IP TTL-check and MUST NOT update the ADspec Break bit. This is because the whole PCN-domain is effectively handled by E2E RSVP as a virtual link on which integrated service is indeed supported (and admission control performed) so that the Break bit MUST NOT be set, see also [draft-lefaucheur-rsvp-ecn-01].

The Deaggregator forwards the E2E Path message towards the receiver.

3.4. Initiation of new Aggregate Path Message by Aggregating Router

To comply with this specification, for the initiation of the new RSVP generic aggregated Path message by the Aggregator (also PCN-ingress-node in this document), the same methods MUST be used as the ones described in [RFC4860].

3.5. Handling Of Aggregate Path Message By Interior Routers

The Aggregate Path messages traverse zero or more PCN-interior-nodes. The PCN-interior-nodes receive the Aggregated Path message on an interior interface and forward it on another interior interface. It is considered that, by configuration, the PCN-interior-nodes ignore the Aggregated Path signaling messages. Therefore, the Aggregated Path messages are simply forwarded as normal IP datagrams.

3.6. Handling Of Aggregate Path Message By Deaggregating Router

When receiving the Aggregated Path message, the Deaggregator (also PCN-egress-node and Decision Point in this document) performs the regular [RFC4860] procedures, augmented with the following rules:

- o) When the received Aggregated Path message by the Deaggregator contains the RSVP-AGGREGATE-IPv4-PCN-response or RSVP-AGGREGATE-IPv6-PCN-response PCN objects, which carry the PCN-sent-rate, then the procedures specified in Section 3.18 of this document MUST be followed.

3.7. Handling of E2E Resv Message by Deaggregating Router

When the E2E Resv message arrives at the exterior interface of the Deaggregator, (also PCN-egress-node and Decision Point in this document) then standard RSVP aggregation [RFC4860] procedures are used, augmented with the following rules:

- o) The E2E RSVP session associated with an E2E Resv message that arrives at the external interface of the Deaggregator is mapped/matched with an RSVP generic aggregate and with a PCN ingress-egress-aggregate.
- o) Depending on the type of the PCN edge behavior supported by the Deaggregator, the PCN admission control procedures specified in Section 3.3.1 of [RFC6661] or [RFC6662] MUST be followed. Since no admission control procedures over the RSVP aggregated reservations in the PCN-core are required, unlike [RFC4860], the Deaggregator does not perform any admission control of the E2E Reservation over the mapped generic aggregate RSVP reservation. If the PCN based admission control procedure is successful then the Deaggregator MUST allow the new flow to be admitted onto the associated RSVP generic aggregation reservation and onto the PCN ingress-egress-aggregate, see [RFC6661] and [RFC6662]. If the PCN based admission control procedure is not successful, then the E2E Resv MUST NOT be admitted onto the associated RSVP generic aggregate reservation and onto the PCN ingress-egress-aggregation. The E2E Resv message is further processed according to [RFC4860].

The way of how the PCN-admission-state is maintained is specified in [RFC6661] and [RFC6662].

3.8. Handling Of E2E Resv Message By Interior Routers

The E2E Resv messages traversing the PCN core are IP addressed to the Aggregating router and are not marked with Router Alert, therefore the E2E Resv messages are simply forwarded as normal IP datagrams.

3.9. Initiation of New Aggregate Resv Message By Deaggregating Router

To comply with this specification, for the initiation of the new RSVP generic aggregated Resv message by the Deaggregator (also PCN-egress-node and Decision Point in this document), the same methods MUST be used as the ones described in Section 4 of [RFC4860] augmented with the following rules:

- o) The size of the generic aggregate reservation is irrelevant, see Section 2.6, and can be set to any arbitrary value by the PCN-egress node. The Deaggregator SHOULD set the value of a RSVP generic aggregate reservation to a null bandwidth. We also observe that there is no need for dynamic adjustment of the RSVP generic aggregated reservation size.

- o) When [RFC6661] is used and the ETM-rate measured by the Deaggregator contains a non-zero value for some ingress-egress-aggregate, see [RFC6661] and [RFC6662], the Deaggregator MUST request the PCN-ingress-node to provide an estimate of the rate (PCN-sent-rate) at which the Aggregator (also PCN-ingress-node in this document) is receiving PCN-traffic that is destined for the given ingress-egress-aggregate.
- o) When [RFC6662] is used and the PCN-admission-state computed by the Deaggregator, on the basis of the CLE is "block" for the given ingress-egress-aggregate, the Deaggregator MUST request the PCN-ingress-node to provide an estimate of the rate (PCN-sent-rate) at which the Aggregator is receiving PCN-traffic that is destined for the given ingress-egress-aggregate.
- o) In the above two cases and when the PCN-sent-rate needs to be requested from the Aggregator, the Deaggregator MUST generate and send an (refresh) Aggregated Resv message to the Aggregator that MUST carry one of the following PCN objects, see Section 4.1, depending on whether IPv4 or IPv6 is supported:
 - o) RSVP-AGGREGATE-IPv4-PCN-request
 - o) RSVP-AGGREGATE-IPv6-PCN-request.

3.10. Handling of Aggregate Resv Message by Interior Routers

The Aggregated Resv messages traversing the PCN core are IP addressed to the Aggregating router and are not marked with Router Alert, therefore the Aggregated Resv messages are simply forwarded as normal IP datagrams.

3.11. Handling of E2E Resv Message by Aggregating Router

When the E2E Resv message arrives at the interior interface of the Aggregator (also PCN-ingress-node in this document), then standard RSVP aggregation [RFC4860] procedures are used.

3.12. Handling of Aggregated Resv Message by Aggregating Router

When the Aggregated Resv message arrives at the interior interface of the Aggregator, (also PCN-ingress-node in this document), then standard RSVP aggregation [RFC4860] procedures are used, augmented with the following rules:

- o) the Aggregator SHOULD use the information carried by the PCN objects, see Section 4, and follow the steps specified in [RFC6661], [RFC6662]. If the "R" flag carried by the RSVP-AGGREGATE-IPv4-PCN-request or RSVP-AGGREGATE-IPv6-PCN-request PCN objects is set to ON, see Section 4.1, then the Aggregator follows the steps described in Section 3.4 of [RFC6661] and [RFC6662] on calculating the PCN-sent-rate. In particular, the Aggregator MUST provide the estimated current rate of PCN-traffic received at that node and destined for a given ingress-egress-aggregate in octets per second (the PCN-sent-rate). The way this rate estimate is derived is a matter of implementation, see [RFC6661] or [RFC6662].

- o) the Aggregator initiates an Aggregated Path message. In particular, when the Aggregator receives an Aggregated Resv message which carries one of the following PCN objects: RSVP-AGGREGATE-IPv4-PCN-request or RSVP-AGGREGATE-IPv6-PCN-request, with the flag "R" set to ON, see Section 4.1, the Aggregator initiates an Aggregated Path message, and includes the calculated PCN-sent-rate into the RSVP-AGGREGATE-IPv4-PCN-response or RSVP-AGGREGATE-IPv6-PCN-response PCN objects, see Section 4.1, which that MUST be carried by the Aggregated Path message. This Aggregated Path message is sent towards the Deaggregator (also PCN-egress-node and Decision Point in this document) that requested the calculation of the PCN-sent-rate.

3.13. Removal of E2E Reservation

To comply with this specification, for the removal of E2E reservations, the same methods MUST be used as the ones described in Section 4 of [RFC4860] and [RFC4495].

3.14. Removal of Aggregate Reservation

To comply with this specification, for the removal of RSVP generic aggregated reservations, the same methods MUST be used as the ones described in Section 4 of [RFC4860] and Section 2.10 of [RFC3175]. In particular, should an aggregate reservation go away (presumably due to a configuration change, route change, or policy event), the E2E reservations it supports are no longer active. They MUST be treated accordingly.

3.15. Handling of Data On Reserved E2E Flow by Aggregating Router

The handling of data on the reserved E2E flow by Aggregator (also PCN-ingress-node in this document) uses the procedures described in [RFC4860] augmented with:

- o) Regarding, PCN marking and traffic classification the procedures defined in Section 2.2 and 2.3 of this document are used.

3.16. Procedures for Multicast Sessions

In this document no multicast sessions are considered.

3.17. Misconfiguration of PCN-node

In an event where a PCN-node is misconfigured within a PCN-domain, the desired behavior is same as described in Section 3.10.

3.18 PCN based Flow Termination

When the Deaggregator (also PCN-egress-node and Decision Point in this document) needs to terminate an amount of traffic associated with one ingress-egress-aggregate (see Section 3.3.2 of [RFC6661] and [RFC6662]), then several procedures of terminating E2E microflows can be deployed. The default procedure of terminating E2E microflows (i.e., PCN-flows) is as follows, see i.e., [RFC6661] and [RFC6662].

For the same ingress-egress-aggregate, select a number of E2E microflows to be terminated in order to decrease the total incoming amount of bandwidth associated with one ingress-egress-aggregate by the amount of traffic to be terminated, see above. In this situation the same mechanisms for terminating an E2E microflow can be followed as specified in [RFC2205]. However, based on a local policy, the Deaggregator could use other ways of selecting which microflows should be terminated. For example, for the same ingress-egress-aggregate, select a number of E2E microflows to be terminated or to reduce their reserved bandwidth in order to decrease the total incoming amount of bandwidth associated with one ingress-egress-aggregate by the amount of traffic to be terminated. In this situation the same mechanisms for terminating an E2E microflow or reducing bandwidth associated with an E2E microflow can be followed as specified in [RFC4495].

4. Protocol Elements

The protocol elements in this document are using the ones defined in Section 4 of [RFC4860] and Section 3 of [RFC3175] augmented with the following rules:

- o) the DSCP value included in the SESSION object, SHOULD be set equal to a PCN-compatible Diffserv codepoint.
- o) Extended vDstPort SHOULD be set to the IPv4 or IPv6 destination addresses, of the Aggregator (also PCN-ingress-node in this document), see [RFC4860].
- o) When the Deaggregator (also PCN-egress-node and Decision Point in this document) needs to request the PCN-sent-rate from the PCN-ingress-node, see Section 3.9 of this document, the Deaggregator MUST generate and send an (refresh) Aggregate Resv message to the Aggregator that MUST carry one of the following PCN objects, see Section 4.1, depending on whether IPv4 or IPv6 is supported:
 - o) RSVP-AGGREGATE-IPv4-PCN-request
 - o) RSVP-AGGREGATE-IPv6-PCN-request.
- o) When the Aggregator receives an Aggregate Resv message which carries one of the following PCN objects:
RSVP-AGGREGATE-IPv4-PCN-request or
RSVP-AGGREGATE-IPv6-PCN-request, with the flag "R" set to ON, see Section 4.1, then the Aggregator MUST generate and send to the Deaggregator an Aggregated Path message which carries one of the following PCN objects, see Section 4.1, depending on whether IPv4 or IPv6 is supported:
 - o) RSVP-AGGREGATE-IPv4-PCN-response,
 - o) RSVP-AGGREGATE-IPv6-PCN-response.

4.1 PCN objects

This section describes four types of PCN objects that can be carried by the (refresh) Aggregate Path or the (refresh) Aggregate Resv messages specified in [RFC4860].

These objects are:

- o RSVP-AGGREGATE-IPv4-PCN-request,
- o RSVP-AGGREGATE-IPv6-PCN-request,
- o RSVP-AGGREGATE-IPv4-PCN-response,
- o RSVP-AGGREGATE-IPv6-PCN-response.

- o) RSVP-AGGREGATE-IPv4-PCN-request: PCN request object, when IPv4 addresses are used:

Class = 248 (PCN)

C-Type = 1 (RSVP-AGGREGATE-IPv4-PCN-request

+-----+-----+-----+-----+	
	IPv4 PCN-ingress-node Address (4 bytes)
+-----+-----+-----+-----+	
	IPv4 PCN-egress-node Address (4 bytes)
+-----+-----+-----+-----+	
	IPv4 Decision Point Address (4 bytes)
+-----+-----+-----+-----+	
R	Reserved
+-----+-----+-----+-----+	

- o) RSVP-AGGREGATE-IPv6-PCN-request: PCN object, when IPv6 addresses are used:

Class = 248 (PCN)

C-Type = 2 (RSVP-AGGREGATE-IPv6-PCN-request

+-----+-----+-----+-----+	
	IPv6 PCN-ingress-node Address (16 bytes)
+	
+	
+-----+-----+-----+-----+	
	IPv6 PCN-egress-node Address (16 bytes)
+	
+	
+-----+-----+-----+-----+	
	Decision Point Address (16 bytes)
+	
+	
+-----+-----+-----+-----+	
R	Reserved
+-----+-----+-----+-----+	

- o) RSVP-AGGREGATE-IPv4-PCN-response: PCN object, IPv4 addresses are used:
 Class = 248 (PCN)
 C-Type = 3 (RSVP-AGGREGATE-IPv4-PCN-response)

```

+-----+-----+-----+-----+
| IPv4 PCN-ingress-node Address (4 bytes) |
+-----+-----+-----+-----+
| IPv4 PCN-egress-node Address (4 bytes) |
+-----+-----+-----+-----+
| IPv4 Decision Point Address (4 bytes) |
+-----+-----+-----+-----+
| PCN-sent-rate |
+-----+-----+-----+-----+

```

- o) RSVP-AGGREGATE-IPv6-PCN-response: PCN object, IPv6 addresses are used:
 Class = 248 (PCN)
 C-Type = 4 (RSVP-AGGREGATE-IPv6-PCN-response)

```

+-----+-----+-----+-----+
|                                     |
+                                     +
| IPv6 PCN-ingress-node Address (16 bytes) |
+                                     +
|                                     |
+-----+-----+-----+-----+
|                                     |
+                                     +
| IPv6 PCN-egress-node Address (16 bytes) |
+                                     +
|                                     |
+-----+-----+-----+-----+
|                                     |
+                                     +
| Decision Point Address (16 bytes) |
+                                     +
|                                     |
+-----+-----+-----+-----+
| PCN-sent-rate |
+-----+-----+-----+-----+

```

The fields carried by the PCN object are specified in [RFC6663], [RFC6661] and [RFC6662]:

- o the IPv4 or IPv6 address of the PCN-ingress-node (Aggregator) and the IPv4 or IPv6 address of the PCN-egress-node (Deaggregator); together they specify the ingress-egress-aggregate to which the report refers. According to [RFC6663] the report should carry the identifier of the PCN-ingress-node (Aggregator) and the identifier of the PCN-egress-node (Deaggregator) (typically their IP addresses);
- o Decision Point address specify the IPv4 or IPv6 address of the Decision Point. In this document this field MUST contain the IP address of the Deaggregator.
- o "R": 1 bit flag that when set to ON, signifies, according to [RFC6661] and [RFC6662], that the PCN-ingress-node (Aggregator) MUST provide an estimate of the rate (PCN-sent-rate) at which the PCN-ingress-node (Aggregator) is receiving PCN-traffic that is destined for the given ingress-egress-aggregate.
- o "Reserved": 31 bits that are currently not used by this document and are reserved. These SHALL be set to 0 and SHALL be ignored on reception.
- o PCN-sent-rate: the PCN-sent-rate for the given ingress-egress-aggregate. It is expressed in octets/second; its format is a 32-bit IEEE floating point number; The PCN-sent-rate is specified in [RFC6661] and [RFC6662] and it represents the estimate of the rate at which the PCN-ingress-node (Aggregator) is receiving PCN-traffic that is destined for the given ingress-egress-aggregate.

5. Security Considerations

The security considerations specified in [RFC2205], [RFC4860] and [RFC5559] apply to this document. In addition, [RFC4230] and [RFC6411] provide useful guidance on RSVP security mechanisms.

Security within a PCN domain is fundamentally based on the controlled environment trust assumption stated in Section 6.3.1 of [RFC5559], in particular that all PCN-nodes are PCN-enabled and are trusted to perform accurate PCN-metering and PCN-marking.

In the PCN domain environments addressed by this document, Generic Aggregate Resource ReSerVation Protocol (RSVP) messages specified in [RFC4860] are used for support of the PCN Controlled Load (CL) and Single Marking (SM) edge behaviors over a Diffserv cloud using Pre-Congestion Notification. Hence the security mechanisms discussed in [RFC4860] are applicable. Specifically, the INTEGRITY object [RFC2747][RFC3097] can be used to provide hop-by-hop RSVP message integrity, node authentication and replay protection, thereby protecting against corruption and spoofing of RSVP messages and PCN feedback conveyed by RSVP messages.

For these reasons, this document does not introduce significant additional security considerations beyond those discussed in

[RFC5559] and [RFC4860].

6. IANA Considerations

IANA has modified the RSVP parameters registry, 'Class Names, Class Numbers, and Class Types' subregistry, to add a new Class Number and assign 4 new C-Types under this new Class Number, as described below, see Section 4.1:

Class Number	Class Name	Reference
-----	-----	-----
248	PCN	this document
Class Types or C-Types:		
1	RSVP-AGGREGATE-IPv4-PCN-request	this document
2	RSVP-AGGREGATE-IPv6-PCN-request	this document
3	RSVP-AGGREGATE-IPv4-PCN-response	this document
4	RSVP-AGGREGATE-IPv6-PCN-response	this document

When this draft is published as an RFC, IANA should update the reference for the above 5 items to that published RFC (and the RFC Editor should remove this sentence).

7. Acknowledgments

We would like to thank the authors of [draft-lefaucheur-rsvp-ecn-01.txt], since some ideas used in this document are based on the work initiated in [draft-lefaucheur-rsvp-ecn-01.txt]. Moreover, we would like to thank Bob Briscoe, David Black, Ken Carlberg, Tom Taylor, Philip Eardley, Michael Menth, Toby Moncaster, James Polk, Scott Bradner, Lixia Zhang and Robert Sparks for the provided comments. In particular, we would like to thank Francois Le Faucheur for contributing in addition to comments also to a significant amount of text.

8. Normative References

- [RFC6661] T. Taylor, A. Charny, F. Huang, G. Karagiannis, M. Menth, "PCN Boundary Node Behaviour for the Controlled Load (CL) Mode of Operation", July 2012.
- [RFC6662] A. Charny, J. Zhang, G. Karagiannis, M. Menth, T. Taylor, "PCN Boundary Node Behaviour for the Single Marking (SM) Mode of Operation", July 2012.
- [RFC6663] G. Karagiannis, T. Taylor, K. Chan, M. Menth, P. Eardley, " Requirements for Signaling of (Pre-) Congestion Information in a DiffServ Domain", July 2012.

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2205] Braden, R., ed., et al., "Resource ReSerVation Protocol (RSVP)- Functional Specification", RFC 2205, September 1997.
- [RFC3140] Black, D., Brim, S., Carpenter, B., and F. Le Faucheur, "Per Hop Behavior Identification Codes", RFC 3140, June 2001.
- [RFC3175] Baker, F., Iturralde, C., Le Faucheur, F., and B. Davie, "Aggregation of RSVP for IPv4 and IPv6 Reservations", RFC 3175, September 2001.
- [RFC4495] Polk, J. and S. Dhesikan, "A Resource Reservation Protocol (RSVP) Extension for the Reduction of Bandwidth of a Reservation Flow", RFC 4495, May 2006.
- [RFC4860] F. Le Faucheur, B. Davie, P. Bose, C. Christou, M. Davenport, "Generic Aggregate Resource ReSerVation Protocol (RSVP) Reservations", RFC4860, May 2007.
- [RFC5670] Eardley, P., "Metering and Marking Behaviour of PCN-Nodes", RFC 5670, November 2009.
- [RFC6660] Moncaster, T., Briscoe, B., and M. Menth, "Baseline Encoding and Transport of Pre-Congestion Information", RFC 6660, July 2012.

9. Informative References

- [draft-lefaucheur-rsvp-ecn-01.txt] Le Faucheur, F., Charny, A., Briscoe, B., Eardley, P., Chan, K., and J. Babiarz, "RSVP Extensions for Admission Control over Diffserv using Pre-congestion Notification (PCN) (Work in progress)", June 2006.
- [RFC1633] Braden, R., Clark, D., and S. Shenker, "Integrated Services in the Internet Architecture: an Overview", RFC 1633, June 1994.
- [RFC2211] J. Wroclawski, Specification of the Controlled-Load Network Element Service, September 1997
- [RFC2212] S. Shenker et al., Specification of Guaranteed Quality of Service, September 1997
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, December 1998.
- [RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z. and W. Weiss, "A framework for Differentiated Services", RFC 2475, December 1998.

[RFC2747] Baker, F., Lindell, B., and M. Talwar, "RSVP Cryptographic Authentication", RFC 2747, January 2000.

[RFC2753] Yavatkar, R., D. Pendarakis and R. Guerin, "A Framework for Policy-based Admission Control", January 2000.

[RFC2998] Bernet, Y., Yavatkar, R., Ford, P., Baker, F., Zhang, L., Speer, M., Braden, R., Davie, B., Wroclawski, J. and E. Felstaine, "A Framework for Integrated Services Operation Over DiffServ Networks", RFC 2998, November 2000.

[RFC3097] Braden, R. and L. Zhang, "RSVP Cryptographic Authentication -- Updated Message Type Value", RFC 3097, April 2001.

[RFC4230] H. Tschofenig, R. Graveman, "RSVP Security Properties", RFC 4230, December 2005.

[RFC5559] Eardley, P., "Pre-Congestion Notification (PCN) Architecture", RFC 5559, June 2009.

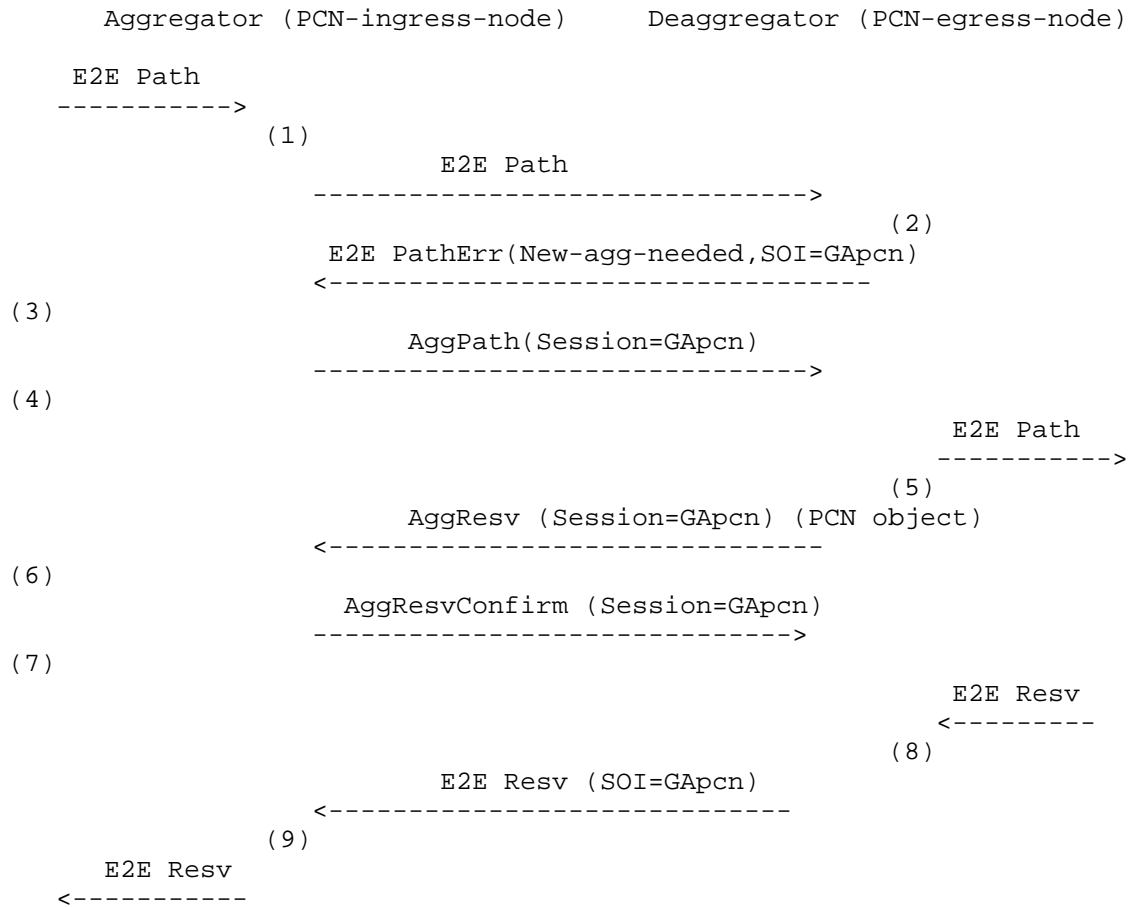
[RFC6411] M. Behringer, F. Le Faucheur, B. Weis, "Applicability of Keying Methods for RSVP Security", RFC 6411, October 2011.

[SIG-NESTED] Baker, F. and P. Bose, "QoS Signaling in a Nested Virtual Private Network", Work in Progress, July 2007.

10. Appendix A: Example Signaling Flow

This appendix is based on the appendix provided in [RFC4860]. In particular, it provides an example signaling flow of the specification detailed in Section 3 and 4.

This signaling flow assumes an environment where E2E reservations are aggregated over generic aggregate RSVP reservations and applied over a PCN domain. In particular the Aggregator (PCN-ingress-node) and Deaggregator (PCN-egress-node) are located at the boundaries of the PCN domain. The PCN-interior-nodes are located within the PCN-domain, between the PCN-boundary nodes, but are not shown in this Figure. It illustrates a possible RSVP message flow that could take place in the successful establishment of a unicast E2E reservation that is the first between a given pair of Aggregator/Deaggregator.



(1) The Aggregator forwards E2E Path into the aggregation region after modifying its IP protocol number to RSVP-E2E-IGNORE

(2) Let's assume no Aggregate Path exists. To be able to accurately update the ADSPEC of the E2E Path, the Deaggregator needs the ADSPEC of Aggregate Path. In this example, the Deaggregator elects to instruct the Aggregator to set up an Aggregate Path state for the PCN PHB-ID. To do that, the Deaggregator sends an E2E PathErr message with a New-Agg-Needed PathErr code.

The PathErr message also contains a SESSION-OF-INTEREST (SOI) object. The SOI contains a GENERIC-AGGREGATE SESSION (GApcn) whose PHB-ID is set to the PCN PHB-ID. The GENERIC-AGGREGATE SESSION contains an interface-independent Deaggregator address inside the DestAddress and appropriate values inside the vDstPort and Extended vDstPort fields. In this document, the Extended vDstPort SHOULD contain the IPv4 or IPv6 address of the Aggregator.

(3) The Aggregator follows the request from the Deaggregator and

signals an Aggregate Path for the GENERIC-AGGREGATE Session (GApcn).

- (4) The Deaggregator takes into account the information contained in the ADSPEC from both Aggregate Paths and updates the E2E Path ADSPEC accordingly. The PCN-egress-node MUST NOT perform the RSVP-TTL vs IP TTL-check and MUST NOT update the ADSpec Break bit. This is because the whole PCN-domain is effectively handled by E2E RSVP as a virtual link on which integrated service is indeed supported (and admission control performed) so that the Break bit MUST NOT be set, see also [draft-lefaucheur-rsvp-ecn-01]. The Deaggregator also modifies the E2E Path IP protocol number to RSVP before forwarding it.
- (5) In this example, the Deaggregator elects to immediately proceed with establishment of the generic aggregate reservation. In effect, the Deaggregator can be seen as anticipating the actual demand of E2E reservations so that the generic aggregate reservation is in place when the E2E Resv request arrives, in order to speed up establishment of E2E reservations. Here it is also assumed that the Deaggregator includes the optional Resv Confirm Request in the Aggregate Resv message.
- (6) The Aggregator merely complies with the received ResvConfirm Request and returns the corresponding Aggregate ResvConfirm.
- (7) The Deaggregator has explicit confirmation that the generic aggregate reservation is established.
- (8) On receipt of the E2E Resv, the Deaggregator applies the mapping policy defined by the network administrator to map the E2E Resv onto a generic aggregate reservation. Let's assume that this policy is such that the E2E reservation is to be mapped onto the generic aggregate reservation with the PCN PHB-ID=x. The Deaggregator knows that a generic aggregate reservation (GApcn) is in place for the corresponding PHB-ID since (7). At this step the Deaggregator maps the generic aggregated reservation onto one ingress-egress-aggregate maintained by the Deaggregator (as a PCN-egress-node), see Section 3.7. The Deaggregator performs admission control of the E2E Resv onto the generic Aggregate reservation for the PCN PHB-ID (GApcn). The Deaggregator takes also into account the PCN admission control procedure as specified in [RFC6661] and [RFC6662], see Section 3.7. If one or both the admission control procedures (PCN based admission control procedure and admission control procedure specified in [RFC4860]) are not successful, then the E2E Resv is not admitted onto the associated RSVP generic aggregate reservation for the PCN PHB-ID (GApcn). Otherwise, assuming that the generic aggregate reservation for the PCN (GApcn) had been established with sufficient bandwidth to support the E2E Resv, the Deaggregator adjusts its counter, tracking the unused bandwidth on the generic aggregate reservation. Then it forwards the E2E Resv to the Aggregator including a SESSION-OF-INTEREST

object conveying the selected mapping onto GApn (and hence onto the PCN PHB-ID).

- (9) The Aggregator records the mapping of the E2E Resv onto GApn (and onto the PCN PHB-ID). The Aggregator removes the SOI object and forwards the E2E Resv towards the sender.

11. Authors' Address

Georgios Karagiannis
Huawei Technologies
Hansaallee 205,
40549 Dusseldorf,
Germany
Email: Georgios.Karagiannis@huawei.com

Anurag Bhargava
Cisco Systems, Inc.
7100-9 Kit Creek Road
PO Box 14987
RESEARCH TRIANGLE PARK, NORTH CAROLINA 27709-4987
USA
Email: anuragb@cisco.com

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: July 28, 2015

M. Tuexen
Muenster Univ. of Appl. Sciences
R. Stewart
Netflix, Inc.
R. Jesup
WorldGate Communications
S. Loreto
Ericsson
January 24, 2015

DTLS Encapsulation of SCTP Packets
draft-ietf-tsvwg-sctp-dtls-encaps-09.txt

Abstract

The Stream Control Transmission Protocol (SCTP) is a transport protocol originally defined to run on top of the network protocols IPv4 or IPv6. This document specifies how SCTP can be used on top of the Datagram Transport Layer Security (DTLS) protocol. Using the encapsulation method described in this document, SCTP is unaware of the protocols being used below DTLS; hence explicit IP addresses cannot be used in the SCTP control chunks. As a consequence, the SCTP associations carried over DTLS can only be single homed.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on July 28, 2015.

Copyright Notice

Copyright (c) 2015 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Overview	2
2. Conventions	3
3. Encapsulation and Decapsulation Procedure	3
4. General Considerations	3
5. DTLS Considerations	4
6. SCTP Considerations	5
7. IANA Considerations	6
8. Security Considerations	6
9. Acknowledgments	7
10. References	7
Appendix A. NOTE to the RFC-Editor	9
Authors' Addresses	9

1. Overview

The Stream Control Transmission Protocol (SCTP) as defined in [RFC4960] is a transport protocol running on top of the network protocols IPv4 [RFC0791] or IPv6 [RFC2460]. This document specifies how SCTP is used on top of the Datagram Transport Layer Security (DTLS) protocol. DTLS 1.0 is defined in [RFC4347] and the latest version when this RFC was published, DTLS 1.2, is defined in [RFC6347]. This encapsulation is used for example within the WebRTC protocol suite (see [I-D.ietf-rtcweb-overview] for an overview) for transporting non-SRTP data between browsers. The architecture of this stack is described in [I-D.ietf-rtcweb-data-channel].

[NOTE to RFC-Editor:

Please ensure that the authors double check the above statement about DTLS 1.2 during AUTH48 and then remove this note before publication.

]

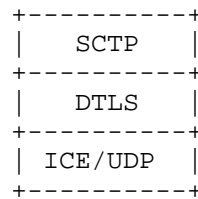


Figure 1: Basic stack diagram

This encapsulation of SCTP over DTLS over UDP or ICE/UDP (see [RFC5245]) can provide a NAT traversal solution in addition to confidentiality, source authentication, and integrity protected transfers. Please note that using ICE does not necessarily imply that a different packet format is used on the wire.

Please note that the procedures defined in [RFC6951] for dealing with the UDP port numbers do not apply here. When using the encapsulation defined in this document, SCTP is unaware about the protocols used below DTLS.

2. 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. Encapsulation and Decapsulation Procedure

When an SCTP packet is provided to the DTLS layer, the complete SCTP packet, consisting of the SCTP common header and a number of SCTP chunks, is handled as the payload of the application layer protocol of DTLS. When the DTLS layer has processed a DTLS record containing a message of the application layer protocol, the payload is passed to the SCTP layer. The SCTP layer expects an SCTP common header followed by a number of SCTP chunks.

4. General Considerations

An implementation of SCTP over DTLS MUST implement and use a path maximum transmission unit (MTU) discovery method that functions without ICMP to provide SCTP/DTLS with an MTU estimate. An implementation of "Packetization Layer Path MTU Discovery" [RFC4821] either in SCTP or DTLS is RECOMMENDED.

The path MTU discovery is performed by SCTP when SCTP over DTLS is used for data channels (see Section 5 of [I-D.ietf-rtcweb-data-channel]).

5. DTLS Considerations

The DTLS implementation MUST support DTLS 1.0 [RFC4347] and SHOULD support the most recently published version of DTLS, which was DTLS 1.2 [RFC6347] when this RFC was published. In the absence of a revision to this document, the latter requirement applies to all future versions of DTLS when they are published as RFCs. This document will only be revised if a revision to DTLS or SCTP makes a revision to the encapsulation necessary.

[NOTE to RFC-Editor:

Please ensure that the authors double check the above statement about DTLS 1.2 during AUTH48 and then remove this note before publication.

]

SCTP performs segmentation and reassembly based on the path MTU. Therefore the DTLS layer MUST NOT use any compression algorithm.

The DTLS MUST support sending messages larger than the current path MTU. This might result in sending IP level fragmented messages.

If path MTU discovery is performed by the DTLS layer, the method described in [RFC4821] MUST be used. For probe packets, the extension defined in [RFC6520] MUST be used.

If path MTU discovery is performed by the SCTP layer and IPv4 is used as the network layer protocol, the DTLS implementation SHOULD allow the DTLS user to enforce that the corresponding IPv4 packet is sent with the Don't Fragment (DF) bit set. If controlling the DF bit is not possible, for example due to implementation restrictions, a safe value for the path MTU has to be used by the SCTP stack. It is RECOMMENDED that the safe value does not exceed 1200 bytes. Please note that [RFC1122] only requires end hosts to be able to reassemble fragmented IP packets up to 576 bytes in length.

The DTLS implementation SHOULD allow the DTLS user to set the Differentiated services code point (DSCP) used for IP packets being sent (see [RFC2474]). This requires the DTLS implementation to pass the value through and the lower layer to allow setting this value. If the lower layer does not support setting the DSCP, then the DTLS user will end up with the default value used by protocol stack. Please note that only a single DSCP value can be used for all packets belonging to the same SCTP association.

Using explicit congestion notifications (ECN) in SCTP requires the DTLS layer to pass the ECN bits through and its lower layer to expose access to them for sent and received packets (see [RFC3168]). The implementation of DTLS and its lower layer have to provide this support. If this is not possible, for example due to implementation restrictions, ECN can't be used by SCTP.

6. SCTP Considerations

This section describes the usage of the base protocol and the applicability of various SCTP extensions.

6.1. Base Protocol

This document uses SCTP [RFC4960] with the following restrictions, which are required to reflect that the lower layer is DTLS instead of IPv4 and IPv6 and that SCTP does not deal with the IP addresses or the transport protocol used below DTLS:

- o A DTLS connection MUST be established before an SCTP association can be set up.
- o Multiple SCTP associations MAY be multiplexed over a single DTLS connection. The SCTP port numbers are used for multiplexing and demultiplexing the SCTP associations carried over a single DTLS connection.
- o All SCTP associations are single-homed, because DTLS does not expose any address management to its upper layer. Therefore it is RECOMMENDED to set the SCTP parameter `path.max.retrans` to `association.max.retrans`.
- o The INIT and INIT-ACK chunk MUST NOT contain any IPv4 Address or IPv6 Address parameters. The INIT chunk MUST NOT contain the Supported Address Types parameter.
- o The implementation MUST NOT rely on processing ICMP or ICMPv6 packets, since the SCTP layer most likely is unable to access the SCTP common header in the plain text of the packet, which triggered the sending of the ICMP or ICMPv6 packet. This applies in particular to path MTU discovery when performed by SCTP.
- o If the SCTP layer is notified about a path change by its lower layers, SCTP SHOULD retest the Path MTU and reset the congestion state to the initial state. The window-based congestion control method specified in [RFC4960], resets the congestion window and slow start threshold to their initial values.

6.2. Padding Extension

When the SCTP layer performs path MTU discovery as specified in [RFC4821], the padding extension defined in [RFC4820] MUST be supported and used for probe packets (HEARTBEAT chunks bundled with PADDING chunks [RFC4820]).

6.3. Dynamic Address Reconfiguration Extension

If the dynamic address reconfiguration extension defined in [RFC5061] is used, ASCONF chunks MUST use wildcard addresses only.

6.4. SCTP Authentication Extension

The SCTP authentication extension defined in [RFC4895] can be used with DTLS encapsulation, but does not provide any additional benefit.

6.5. Partial Reliability Extension

Partial reliability as defined in [RFC3758] can be used in combination with DTLS encapsulation. It is also possible to use additional PR-SCTP policies, for example the ones defined in [I-D.ietf-tsvwg-sctp-prpolicies].

6.6. Stream Reset Extension

The SCTP stream reset extension defined in [RFC6525] can be used with DTLS encapsulation. It is used to reset SCTP streams and add SCTP streams during the lifetime of the SCTP association.

6.7. Interleaving of Large User Messages

SCTP as defined in [RFC4960] does not support the interleaving of large user messages that need to be fragmented and reassembled by the SCTP layer. The protocol extension defined in [I-D.ietf-tsvwg-sctp-ndata] overcomes this limitation and can be used with DTLS encapsulation.

7. IANA Considerations

This document requires no actions from IANA.

8. Security Considerations

Security considerations for DTLS are specified in [RFC4347] and for SCTP in [RFC4960], [RFC3758], and [RFC6525]. The combination of SCTP and DTLS introduces no new security considerations.

SCTP should not process the IP addresses used for the underlying communication since DTLS provides no guarantees about them.

It should be noted that the inability to process ICMP or ICMPv6 messages does not add any security issue. When SCTP is carried over a connection-less lower layer like IPv4, IPv6, or UDP, processing of these messages is required to protect other nodes not supporting SCTP. Since DTLS provides a connection-oriented lower layer, this kind of protection is not necessary.

9. Acknowledgments

The authors wish to thank David Black, Benoit Claise, Spencer Dawkins, Francis Dupont, Gorrry Fairhurst, Stephen Farrell, Christer Holmberg, Barry Leiba, Eric Rescorla, Tom Taylor, Joe Touch and Magnus Westerlund for their invaluable comments.

10. References

10.1. Normative References

- [RFC1122] Braden, R., "Requirements for Internet Hosts - Communication Layers", STD 3, RFC 1122, October 1989.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4347] Rescorla, E. and N. Modadugu, "Datagram Transport Layer Security", RFC 4347, April 2006.
- [RFC4820] Tuexen, M., Stewart, R., and P. Lei, "Padding Chunk and Parameter for the Stream Control Transmission Protocol (SCTP)", RFC 4820, March 2007.
- [RFC4821] Mathis, M. and J. Heffner, "Packetization Layer Path MTU Discovery", RFC 4821, March 2007.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.
- [RFC6347] Rescorla, E. and N. Modadugu, "Datagram Transport Layer Security Version 1.2", RFC 6347, January 2012.
- [RFC6520] Seggelmann, R., Tuexen, M., and M. Williams, "Transport Layer Security (TLS) and Datagram Transport Layer Security (DTLS) Heartbeat Extension", RFC 6520, February 2012.

10.2. Informative References

- [RFC0791] Postel, J., "Internet Protocol", STD 5, RFC 791, September 1981.
- [RFC2460] Deering, S. and R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification", RFC 2460, December 1998.
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, December 1998.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.
- [RFC3758] Stewart, R., Ramalho, M., Xie, Q., Tuexen, M., and P. Conrad, "Stream Control Transmission Protocol (SCTP) Partial Reliability Extension", RFC 3758, May 2004.
- [RFC4895] Tuexen, M., Stewart, R., Lei, P., and E. Rescorla, "Authenticated Chunks for the Stream Control Transmission Protocol (SCTP)", RFC 4895, August 2007.
- [RFC5061] Stewart, R., Xie, Q., Tuexen, M., Maruyama, S., and M. Kozuka, "Stream Control Transmission Protocol (SCTP) Dynamic Address Reconfiguration", RFC 5061, September 2007.
- [RFC5245] Rosenberg, J., "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols", RFC 5245, April 2010.
- [RFC6525] Stewart, R., Tuexen, M., and P. Lei, "Stream Control Transmission Protocol (SCTP) Stream Reconfiguration", RFC 6525, February 2012.
- [RFC6951] Tuexen, M. and R. Stewart, "UDP Encapsulation of Stream Control Transmission Protocol (SCTP) Packets for End-Host to End-Host Communication", RFC 6951, May 2013.
- [I-D.ietf-rtcweb-overview] Alvestrand, H., "Overview: Real Time Protocols for Browser-based Applications", draft-ietf-rtcweb-overview-13 (work in progress), November 2014.

[I-D.ietf-rtcweb-data-channel]

Jesup, R., Loreto, S., and M. Tuexen, "WebRTC Data Channels", draft-ietf-rtcweb-data-channel-13 (work in progress), January 2015.

[I-D.ietf-tsvwg-sctp-prpolicies]

Tuexen, M., Seggelmann, R., Stewart, R., and S. Loreto, "Additional Policies for the Partial Reliability Extension of the Stream Control Transmission Protocol", draft-ietf-tsvwg-sctp-prpolicies-06 (work in progress), December 2014.

[I-D.ietf-tsvwg-sctp-ndata]

Stewart, R., Tuexen, M., Loreto, S., and R. Seggelmann, "Stream Schedulers and a New Data Chunk for the Stream Control Transmission Protocol", draft-ietf-tsvwg-sctp-ndata-02 (work in progress), January 2015.

Appendix A. NOTE to the RFC-Editor

Although the references to [I-D.ietf-tsvwg-sctp-prpolicies] and [I-D.ietf-tsvwg-sctp-ndata] are informative, put this document in REF-HOLD until these two references have been approved and update these references to the corresponding RFCs.

Authors' Addresses

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

Email: tuexen@fh-muenster.de

Randall R. Stewart
Netflix, Inc.
Chapin, SC 29036
US

Email: randall@lakerest.net

Randell Jesup
WorldGate Communications
3800 Horizon Blvd, Suite #103
Trevose, PA 19053-4947
US

Phone: +1-215-354-5166
Email: randell_ietf@jesup.org

Salvatore Loreto
Ericsson
Hirsalantie 11
Jorvas 02420
FI

Email: Salvatore.Loreto@ericsson.com

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: August 20, 2016

Y. Nishida
GE Global Research
P. Natarajan
Cisco Systems
A. Caro
BBN Technologies
P. Amer
University of Delaware
K. Nielsen
Ericsson
February 17, 2016

SCTP-PF: Quick Failover Algorithm in SCTP
draft-ietf-tsvwg-sctp-failover-16.txt

Abstract

SCTP supports multi-homing. However, when the failover operation specified in RFC4960 is followed, there can be significant delay and performance degradation in the data transfer path failover. To overcome this problem this document specifies a quick failover algorithm (SCTP-PF) based on the introduction of a Potentially Failed (PF) state in SCTP Path Management.

The document also specifies a dormant state operation of SCTP. This dormant state operation is required to be followed by an SCTP-PF implementation, but it may equally well be applied by a standard RFC4960 SCTP implementation.

Additionally, the document introduces an alternative switchback operation mode called Primary Path Switchover that will be beneficial in certain situations. This mode of operation applies to both a standard RFC4960 SCTP implementation as well as to a SCTP-PF implementation.

The procedures defined in the document require only minimal modifications to the RFC4960 specification. The procedures are sender-side only and do not impact the SCTP receiver.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on August 20, 2016.

Copyright Notice

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

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

Table of Contents

1. Introduction	3
2. Conventions and Terminology	4
3. SCTP with Potentially Failed Destination State (SCTP-PF) . .	4
3.1. Overview	4
3.2. Specification of the SCTP-PF Procedures	5
4. Dormant State Operation	9
4.1. SCTP Dormant State Procedure	10
5. Primary Path Switchover	11
6. Suggested SCTP Protocol Parameter Values	12
7. Socket API Considerations	12
7.1. Support for the Potentially Failed Path State	13
7.2. Peer Address Thresholds (SCTP_PEER_ADDR_THLDS) Socket Option	14
7.3. Exposing the Potentially Failed Path State (SCTP_EXPOSE_POTENTIALLY_FAILED_STATE) Socket Option . .	15
8. Security Considerations	15
9. MIB Considerations	16
10. IANA Considerations	16
11. Acknowledgements	16
12. Proposed Change of Status (to be Deleted before Publication)	17
13. References	17

13.1. Normative References	17
13.2. Informative References	17
Appendix A. Discussions of Alternative Approaches	18
A.1. Reduce Path.Max.Retrans (PMR)	18
A.2. Adjust RTO related parameters	19
Appendix B. Discussions for Path Bouncing Effect	20
Appendix C. SCTP-PF for SCTP Single-homed Operation	20
Authors' Addresses	21

1. Introduction

The Stream Control Transmission Protocol (SCTP) specified in [RFC4960] supports multi-homing at the transport layer. SCTP's multi-homing features include failure detection and failover procedures to provide network interface redundancy and improved end-to-end fault tolerance. In SCTP's current failure detection procedure, the sender must experience Path.Max.Retrans (PMR) number of consecutive failed timer-based retransmissions on a destination address before detecting a path failure. Until detecting the path failure, the sender continues to transmit data on the failed path. The prolonged time in which [RFC4960] SCTP continues to use a failed path severely degrades the performance of the protocol. To address this problem, this document specifies a quick failover algorithm (SCTP-PF) based on the introduction of a new Potentially Failed (PF) path state in SCTP path management. The performance deficiencies of the [RFC4960] failover operation, and the improvements obtainable from the introduction of a Potentially Failed state in SCTP, were proposed and documented in [NATARAJAN09] for Concurrent Multipath Transfer SCTP [IYENGAR06].

While SCTP-PF can accelerate failover process and improve performance, the risks that an SCTP endpoint enters the dormant state where all destination addresses are inactive can be increased. [RFC4960] leaves the protocol operation during dormant state to implementations and encourages to avoid entering the state as much as possible by careful tuning of the Path.Max.Retrans (PMR) and Association.Max.Retrans (AMR) parameters. We specify a dormant state operation for SCTP-PF which makes SCTP-PF provide the same disruption tolerance as [RFC4960] despite that the dormant state may be entered more quickly. The dormant state operation may equally well be applied by an [RFC4960] implementation and will here serve to provide added fault tolerance for situations where the tuning of the Path.Max.Retrans (PMR) and Association.Max.Retrans (AMR) parameters fail to provide adequate prevention of the entering of the dormant state.

The operation after the recovery of a failed path also impacts the performance of the protocol. With the procedures specified in

[RFC4960] SCTP will, after a failover from the primary path, switch back to use the primary path for data transfer as soon as this path becomes available again. From a performance perspective such a forced switchback of the data transmission path can be suboptimal as the CWND towards the original primary destination address has to be rebuilt once data transfer resumes, [CARO02]. As an optional alternative to the switchback operation of [RFC4960], this document specifies an alternative Primary Path Switchover procedure which avoid such forced switchbacks of the data transfer path. The Primary Path Switchover operation was originally proposed in [CARO02].

While SCTP-PF primarily is motivated by a desire to improve the multi-homed operation, the feature applies also to SCTP single-homed operation. Here the algorithm serves to provide increased failure detection on idle associations, whereas the failover or switchback aspects of the algorithm will not be activated. This is discussed in more detail in Appendix C.

A brief description of the motivation for the introduction of the Potentially Failed state including a discussion of alternative approaches to mitigate the deficiencies of the [RFC4960] failover operation are given in the Appendices. Discussion of path bouncing effects that might be caused by frequent switchovers, are also provided there.

2. Conventions and Terminology

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

3. SCTP with Potentially Failed Destination State (SCTP-PF)

3.1. Overview

To minimize the performance impact during failover, the sender should avoid transmitting data to a failed destination address as early as possible. In the [RFC4960] SCTP path management scheme, the sender stops transmitting data to a destination address only after the destination address is marked inactive. This process takes a significant amount of time as it requires the error counter of the destination address to exceed the Path.Max.Retrans (PMR) threshold. The issue cannot simply be mitigated by lowering of the PMR threshold because this may result in spurious failure detection and unnecessary prevention of the usage of a preferred primary path. Also due to the coupled tuning of the Path.Max.Retrans (PMR) and the Association.Max.Retrans (AMR) parameter values in [RFC4960], lowering

of the PMR threshold may result in lowering of the AMR threshold, which would result in decrease of the fault tolerance of SCTP.

The solution provided in this document is to extend the SCTP path management scheme of [RFC4960] by the addition of the Potentially Failed (PF) state as an intermediate state in between the active and inactive state of a destination address in the [RFC4960] path management scheme, and let the failover of data transfer away from a destination address be driven by the entering of the PF state instead of by the entering of the inactive state. Thereby SCTP may perform quick failover without negatively impacting the overall fault tolerance of [RFC4960] SCTP. At the same time, RTO-based HEARTBEAT probing is initiated towards a destination address once it enters PF state. Thereby SCTP may quickly ascertain whether network connectivity towards the destination address is broken or whether the failover was spurious. In the case where the failover was spurious data transfer may quickly resume towards the original destination address.

The new failure detection algorithm assumes that loss detected by a timeout implies either severe congestion or network connectivity failure. It recommends that by default a destination address is classified as PF at the occurrence of the first timeout.

3.2. Specification of the SCTP-PF Procedures

The SCTP-PF operation is specified as follows:

1. The sender maintains a new tunable SCTP Protocol Parameter called PotentiallyFailed.Max.Retrans (PFMR). The PFMR defines the new intermediate PF threshold on the destination address error counter. When this threshold is exceeded the destination address is classified as PF. The RECOMMENDED value of PFMR is 0. If PFMR is set to be greater than or equal to Path.Max.Retrans (PMR), the resulting PF threshold will be so high that the destination address will reach the inactive state before it can be classified as PF.
2. The error counter of an active destination address is incremented or cleared as specified in [RFC4960]. This means that the error counter of the destination address in active state will be incremented each time the T3-rtx timer expires, or each time a HEARTBEAT chunk is sent when idle and not acknowledged within an RTO. When the value in the destination address error counter exceeds PFMR, the endpoint MUST mark the destination address as in the PF state.

3. A SCTP-PF sender SHOULD NOT send data to destination addresses in PF state when alternative destination addresses in active state are available. Specifically this means that:
 - i When there is outbound data to send and the destination address presently used for data transmission is in PF state, the sender SHOULD choose a destination address in active state, if one exists, and use this destination address for data transmission.
 - ii As specified in [RFC4960] section 6.4.1, when the sender retransmits data that has timed out, it should attempt to pick a new destination address for data retransmission. In this case, the sender SHOULD choose an alternate destination transport address in active state if one exists.
 - iii When there is outbound data to send and the SCTP user explicitly requests to send data to a destination address in PF state, the sender SHOULD send the data to an alternate destination address in active state if one exists.

When choosing among multiple destination addresses in active state an SCTP sender will follow the guiding principles of section 6.4.1 of [RFC4960] of choosing most divergent source-destination pairs compared with, for i.: the destination address in PF state that it performs a failover from, and for ii.: the destination address towards which the data timed out. Rules for picking the most divergent source-destination pair are an implementation decision and are not specified within this document.

In all cases, the sender MUST NOT change the state of chosen destination address, whether this state be active or PF, and it MUST NOT clear the error counter of the destination address as a result of choosing the destination address for data transmission.

4. When the destination addresses are all in PF state or some in PF state and some in inactive state, the sender MUST choose one destination address in PF state and SHOULD transmit or retransmit data to this destination address using the following rules:
 - A. The sender SHOULD choose the destination in PF state with the lowest error count (fewest consecutive timeouts) for data transmission and transmit or retransmit data to this destination.

- B. When there are multiple destination addresses in PF state with same error count, the sender should let the choice among the multiple destination addresses in PF state with equal error count be based on the [RFC4960], section 6.4.1, principles of choosing most divergent source-destination pairs when executing (potentially consecutive) retransmission. Rules for picking the most divergent source-destination pair are an implementation decision and are not specified within this document.

The sender MUST NOT change the state and the error counter of any destination addresses as the result of the selection.

- 5. The HB.interval of the Path Heartbeat function of [RFC4960] MUST be ignored for destination addresses in PF state. Instead HEARTBEAT chunks are sent to destination addresses in PF state once per RTO. HEARTBEAT chunks SHOULD be sent to destination addresses in PF state, but the sending of HEARTBEATS MUST honor whether the Path Heartbeat function (Section 8.3 of [RFC4960]) is enabled for the destination address or not. I.e., if the Path Heartbeat function is disabled for the destination address in question, HEARTBEATS MUST NOT be sent. Note that when Heartbeat function is disabled, it may take longer to transition a destination address in PF state back to active state.
- 6. HEARTBEATS are sent when a destination address reaches the PF state. When a HEARTBEAT chunk is not acknowledged within the RTO, the sender increments the error counter and exponentially backs off the RTO value. If the error counter is less than PMR, the sender transmits another packet containing the HEARTBEAT chunk immediately after timeout expiration on the previous HEARTBEAT. When data is being transmitted to a destination address in the PF state, the transmission of a HEARTBEAT chunk MAY be omitted in case where the receipt of a SACK of the data or a T3-rtx timer expiration on the data can provide equivalent information, such as the case where the data chunk has been transmitted to a single destination address only. Likewise, the timeout of a HEARTBEAT chunk MAY be ignored if data is outstanding towards the destination address.
- 7. When the sender receives a HEARTBEAT ACK from a HEARTBEAT sent to a destination address in PF state, the sender SHOULD clear the error counter of the destination address and transition the destination address back to active state. However, there may be a situation where HEARTBEAT chunks can go through while DATA chunks cannot. Hence, in a situation where a HEARTBEAT ACK arrives while there is data outstanding towards the destination address to which the HEARTBEAT was sent, then an implementation

MAY choose to not have the HEARTBEAT ACK reset the error counter, but have the error counter reset await the fate of the outstanding data transmission. This situation can happen when data is sent to a destination address in PF state. When the sender resumes data transmission on a destination address after a transition of the destination address from PF to active state, it MUST do this following the prescriptions of Section 7.2 of [RFC4960].

8. Additional (PMR - PFMR) consecutive timeouts on a destination address in PF state confirm the path failure, upon which the destination address transitions to the inactive state. As described in [RFC4960], the sender (i) SHOULD notify the ULP about this state transition, and (ii) transmit HEARTBEAT chunks to the inactive destination address at a lower HB.interval frequency as described in Section 8.3 of [RFC4960] (when the Path Heartbeat function is enabled for the destination address).
9. Acknowledgments for chunks that have been transmitted to multiple destinations (i.e., a chunk which has been retransmitted to a different destination address than the destination address to which the chunk was first transmitted) SHOULD NOT clear the error count for an inactive destination address and SHOULD NOT move a destination address in PF state back to active state, since a sender cannot disambiguate whether the ACK was for the original transmission or the retransmission(s). A SCTP sender MAY clear the error counter and move a destination address back to active state by information other than acknowledgments, when it can uniquely determine which destination, among multiple destination addresses, the chunk reached. This document makes no reference to what such information could consist of, nor how such information could be obtained.
10. Acknowledgments for data chunks that has been transmitted to one destination address only MUST clear the error counter for the destination address and MUST transition a destination address in PF state back to active state. This situation can happen when new data is sent to a destination address in the PF state. It can also happen in situations where the destination address is in the PF state due to the occurrence of a spurious T3-rtx timer and acknowledgments start to arrive for data sent prior to occurrence of the spurious T3-rtx and data has not yet been retransmitted towards other destinations. This document does not specify special handling for detection of or reaction to spurious T3-rtx timeouts, e.g., for special operation vis-a-vis the congestion control handling or data retransmission operation towards a destination address which undergoes a transition from

active to PF to active state due to a spurious T3-rtx timeout. But it is noted that this is an area which would benefit from additional attention, experimentation and specification for single-homed SCTP as well as for multi-homed SCTP protocol operation.

11. When all destination addresses are in inactive state, and SCTP protocol operation thus is said to be in dormant state, the prescriptions given in Section 4 shall be followed.
12. The SCTP stack SHOULD expose the PF state of its destination addresses to the ULP as well as provide the means to notify the ULP of state transitions of its destination addresses from active to PF, and vice-versa. However it is recommended that an SCTP stack implementing SCTP-PF also allows for that the ULP is kept ignorant of the PF state of its destinations and the associated state transitions, thus allowing for retain of the simpler state transition model of RFC4960 in the ULP. For this reason it is recommended that an SCTP stack implementing SCTP-PF also provides the ULP with the means to suppress exposure of the PF state and the associated state transitions.

4. Dormant State Operation

In a situation with complete disruption of the communication in between the SCTP Endpoints, the aggressive HEARTBEAT transmissions of SCTP-PF on destination addresses in PF state may make the association enter dormant state faster than a standard [RFC4960] SCTP implementation given the same setting of Path.Max.Retrans (PMR) and Association.Max.Retrans (AMR). For example, an SCTP association with two destination addresses typically would reach dormant state in half the time of an [RFC4960] SCTP implementation in such situations. This is because a SCTP PF sender will send HEARTBEATS and data retransmissions in parallel with RTO intervals when there are multiple destinations addresses in PF state. This argument presumes that $RTO \ll HB.interval$ of [RFC4960]. With the design goal that SCTP-PF shall provide the same level of disruption tolerance as an [RFC4960] SCTP implementation with the same Path.Max.Retrans (PMR) and Association.Max.Retrans (AMR) setting, we prescribe for that an SCTP-PF implementation SHOULD operate as described below in Section 4.1 during dormant state.

An SCTP-PF implementation MAY choose a different dormant state operation than the one described below in Section 4.1 provided that the solution chosen does not decrease the fault tolerance of the SCTP-PF operation.

The below prescription for SCTP-PF dormant state handling MUST NOT be coupled to the value of the PFMR, but solely to the activation of SCTP-PF logic in an SCTP implementation.

It is noted that the below dormant state operation is considered to provide added disruption tolerance also for an [RFC4960] SCTP implementation, and that it can be sensible for an [RFC4960] SCTP implementation to follow this mode of operation. For an [RFC4960] SCTP implementation the continuation of data transmission during dormant state makes the fault tolerance of SCTP be more robust towards situations where some, or all, alternative paths of an SCTP association approach, or reach, inactive state before the primary path used for data transmission observes trouble.

4.1. SCTP Dormant State Procedure

- a. When the destination addresses are all in inactive state and data is available for transfer, the sender MUST choose one destination and transmit data to this destination address.
- b. The sender MUST NOT change the state of the chosen destination address (it remains in inactive state) and it MUST NOT clear the error counter of the destination address as a result of choosing the destination address for data transmission.
- c. The sender SHOULD choose the destination in inactive state with the lowest error count (fewest consecutive timeouts) for data transmission. When there are multiple destinations with same error count in inactive state, the sender SHOULD attempt to pick the most divergent source - destination pair from the last source - destination pair where failure was observed. Rules for picking the most divergent source-destination pair are an implementation decision and are not specified within this document. To support differentiation of inactive destination addresses based on their error count SCTP will need to allow for increment of the destination address error counters up to some reasonable limit above PMR+1, thus changing the prescriptions of [RFC4960], section 8.3, in this respect. The exact limit to apply is not specified in this document but it is considered reasonable to require for the limit to be an order of magnitude higher than the PMR value. A sender MAY choose to deploy other strategies than the strategy defined here. The strategy to prioritize the last active destination address, i.e., the destination address with the fewest error counts is optimal when some paths are permanently inactive, but suboptimal when a path instability is transient.

5. Primary Path Switchover

The objective of the Primary Path Switchover operation is to allow the SCTP sender to continue data transmission on a new working path even when the old primary destination address becomes active again. This is achieved by having SCTP perform a switchover of the primary path to the new working path if the error counter of the primary path exceeds a certain threshold. This mode of operation can be applied not only to SCTP-PF implementations, but also to [RFC4960] implementations.

The Primary Path Switchover operation requires only sender side changes. The details are:

1. The sender maintains a new tunable parameter, called Primary.Switchover.Max.Retrans (PSMR). For SCTP-PF implementations, the PSMR MUST be set greater or equal to the PFMR value. For [RFC4960] implementations the PSMR MUST be set greater or equal to the PMR value. Implementations MUST reject any other values of PSMR.
2. When the path error counter on a set primary path exceeds PSMR, the SCTP implementation MUST autonomously select and set a new primary path.
3. The primary path selected by the SCTP implementation MUST be the path which at the given time would be chosen for data transfer. A previously failed primary path can be used as data transfer path as per normal path selection when the present data transfer path fails.
4. For SCTP-PF, the recommended value of PSMR is PFMR when Primary Path Switchover operation mode is used. This means that no forced switchback to a previously failed primary path is performed. An SCTP-PF implementation of Primary Path Switchover MUST support the setting of PSMR = PFMR. A SCTP-PF implementation of Primary Path Switchover MAY support setting of PSMR > PFMR.
5. For [RFC4960] SCTP, the recommended value of PSMR is PMR when Primary Path Switchover is used. This means that no forced switchback to a previously failed primary path is performed. A [RFC4960] SCTP implementation of Primary Path Switchover MUST support the setting of PSMR = PMR. An [RFC4960] SCTP implementation of Primary Path Switchover MAY support larger settings of PSMR > PMR.

6. It MUST be possible to disable the Primary Path Switchover operation and obtain the standard switchback operation of [RFC4960].

The manner of switchover operation that is most optimal in a given scenario depends on the relative quality of a set primary path versus the quality of alternative paths available as well as on the extent to which it is desired for the mode of operation to enforce traffic distribution over a number of network paths. I.e., load distribution of traffic from multiple SCTP associations may be sought to be enforced by distribution of the set primary paths with [RFC4960] switchback operation. However as [RFC4960] switchback behavior is suboptimal in certain situations, especially in scenarios where a number of equally good paths are available, an SCTP implementation MAY support also, as alternative behavior, the Primary Path Switchover mode of operation and MAY enable it based on applications' requests.

For an SCTP implementation that implements the Primary Path Switchover operation, this specification RECOMMENDS that the standard RFC4960 switchback operation is retained as the default operation.

6. Suggested SCTP Protocol Parameter Values

This document does not alter the [RFC4960] value recommendation for the SCTP Protocol Parameters defined in [RFC4960].

The following protocol parameter is RECOMMENDED:

PotentiallyFailed.Max.Retrans (PFMR) - 0

7. Socket API Considerations

This section describes how the socket API defined in [RFC6458] is extended to provide a way for the application to control and observe the SCTP-PF behavior as well as the Primary Path Switchover function.

Please note that this section is informational only.

A socket API implementation based on [RFC6458] is, by means of the existing SCTP_PEER_ADDR_CHANGE event, extended to provide the event notification when a peer address enters or leaves the potentially failed state as well as the socket API implementation is extended to expose the potentially failed state of a peer address in the existing SCTP_GET_PEER_ADDR_INFO structure.

Furthermore, two new read/write socket options for the level IPPROTO_SCTP and the name SCTP_PEER_ADDR_THLDS and

SCTP_EXPOSE_POTENTIALLY_FAILED_STATE are defined as described below. The first socket option is used to control the values of the PFMR and PSMP parameters described in Section 3 and in Section 5. The second one controls the exposition of the potentially failed path state.

Support for the SCTP_PEER_ADDR_THLDS and SCTP_EXPOSE_POTENTIALLY_FAILED_STATE socket options need also to be added to the function sctp_opt_info().

7.1. Support for the Potentially Failed Path State

As defined in [RFC6458], the SCTP_PEER_ADDR_CHANGE event is provided if the status of a peer address changes. In addition to the state changes described in [RFC6458], this event is also provided, if a peer address enters or leaves the potentially failed state. The notification as defined in [RFC6458] uses the following structure:

```
struct sctp_paddr_change {
    uint16_t spc_type;
    uint16_t spc_flags;
    uint32_t spc_length;
    struct sockaddr_storage spc_aaddr;
    uint32_t spc_state;
    uint32_t spc_error;
    sctp_assoc_t spc_assoc_id;
}
```

[RFC6458] defines the constants SCTP_ADDR_AVAILABLE, SCTP_ADDR_UNREACHABLE, SCTP_ADDR_REMOVED, SCTP_ADDR_ADDED, and SCTP_ADDR_MADE_PRIM to be provided in the spc_state field. This document defines in addition to that the new constant SCTP_ADDR_POTENTIALLY_FAILED, which is reported if the affected address becomes potentially failed.

The SCTP_GET_PEER_ADDR_INFO socket option defined in [RFC6458] can be used to query the state of a peer address. It uses the following structure:

```
struct sctp_paddrinfo {
    sctp_assoc_t spinfo_assoc_id;
    struct sockaddr_storage spinfo_address;
    int32_t spinfo_state;
    uint32_t spinfo_cwnd;
    uint32_t spinfo_srtt;
    uint32_t spinfo_rto;
    uint32_t spinfo_mtu;
};
```

[RFC6458] defines the constants `SCTP_UNCONFIRMED`, `SCTP_ACTIVE`, and `SCTP_INACTIVE` to be provided in the `spinfo_state` field. This document defines in addition to that the new constant `SCTP_POTENTIALLY_FAILED`, which is reported if the peer address is potentially failed.

7.2. Peer Address Thresholds (`SCTP_PEER_ADDR_THLDS`) Socket Option

Applications can control the SCTP-PF behavior by getting or setting the number of consecutive timeouts before a peer address is considered potentially failed or unreachable. The same socket option is used by applications to set and get the number of timeouts before the primary path is changed automatically by the Primary Path Switchover function. This socket option uses the level `IPPROTO_SCTP` and the name `SCTP_PEER_ADDR_THLDS`.

The following structure is used to access and modify the thresholds:

```
struct sctp_paddrthlds {
    sctp_assoc_t spt_assoc_id;
    struct sockaddr_storage spt_address;
    uint16_t spt_pathmaxrxt;
    uint16_t spt_pathpfthld;
    uint16_t spt_pathcpthld;
};
```

`spt_assoc_id`: This parameter is ignored for one-to-one style sockets. For one-to-many style sockets the application may fill in an association identifier or `SCTP_FUTURE_ASSOC`. It is an error to use `SCTP_{CURRENT|ALL}_ASSOC` in `spt_assoc_id`.

`spt_address`: This specifies which peer address is of interest. If a wild card address is provided, this socket option applies to all current and future peer addresses.

`spt_pathmaxrxt`: Each peer address of interest is considered unreachable, if its path error counter exceeds `spt_pathmaxrxt`.

`spt_pathpfthld`: Each peer address of interest is considered Potentially Failed, if its path error counter exceeds `spt_pathpfthld`.

`spt_pathcpthld`: Each peer address of interest is not considered the primary remote address anymore, if its path error counter exceeds `spt_pathcpthld`. Using a value of `0xffff` disables the selection of a new primary peer address. If an implementation does not support the automatically selection of a new primary address, it should indicate an error with `errno` set to `EINVAL` if a value different

from 0xffff is used in `spt_pathcpthld`. For SCTP-PF, the setting of `spt_pathcpthld < spt_pathpfthld` should be rejected with `errno` set to `EINVAL`. For [RFC4960] SCTP, the setting of `spt_pathcpthld < spt_pathmaxrxt` should be rejected with `errno` set to `EINVAL`. A SCTP-PF implementation may support only setting of `spt_pathcpthld = spt_pathpfthld` and `spt_pathcpthld = 0xffff` and a [RFC4960] SCTP implementation may support only setting of `spt_pathcpthld = spt_pathmaxrxt` and `spt_pathcpthld = 0xffff`. In these cases SCTP shall reject setting of other values with `errno` set to `EINVAL`.

7.3. Exposing the Potentially Failed Path State (`SCTP_EXPOSE_POTENTIALLY_FAILED_STATE`) Socket Option

Applications can control the exposure of the potentially failed path state in the `SCTP_PEER_ADDR_CHANGE` event and the `SCTP_GET_PEER_ADDR_INFO` as described in Section 7.1. The default value is implementation specific.

This socket option uses the level `IPPROTO_SCTP` and the name `SCTP_EXPOSE_POTENTIALLY_FAILED_STATE`.

The following structure is used to control the exposition of the potentially failed path state:

```
struct sctp_assoc_value {
    sctp_assoc_t assoc_id;
    uint32_t assoc_value;
};
```

`assoc_id`: This parameter is ignored for one-to-one style sockets. For one-to-many style sockets the application may fill in an association identifier or `SCTP_FUTURE_ASSOC`. It is an error to use `SCTP_{CURRENT|ALL}_ASSOC` in `assoc_id`.

`assoc_value`: The potentially failed path state is exposed if and only if this parameter is non-zero.

8. Security Considerations

Security considerations for the use of SCTP and its APIs are discussed in [RFC4960] and [RFC6458].

The logic introduced by this document does not impact existing SCTP messages on the wire. Also, this document does not introduce any new SCTP messages on the wire that require new security considerations.

SCTP-PF makes SCTP not only more robust during primary path failure/congestion but also more vulnerable to network connectivity/

congestion attacks on the primary path. SCTP-PF makes it easier for an attacker to trick SCTP to change data transfer path, since the duration of time that an attacker needs to negatively influence the network connectivity is much shorter than [RFC4960]. However, SCTP-PF does not constitute a significant change in the duration of time and effort an attacker needs to keep SCTP away from the primary path. With the standard switchback operation [RFC4960] SCTP resumes data transfer on its primary path as soon as the next HEARTBEAT succeeds.

On the other hand, usage of the Primary Path Switchover mechanism, does change the threat analysis. This is because on-path attackers can force a permanent change of the data transfer path by blocking the primary path until the switchover of the primary path is triggered by the Primary Path Switchover algorithm. This especially will be the case when the Primary Path Switchover is used together with SCTP-PF with the particular setting of PSMR = PFMR = 0, as Primary Path Switchover here happens already at the first RTO timeout experienced. Users of the Primary Path Switchover mechanism should be aware of this fact.

The event notification of path state transfer from active to potentially failed state and vice versa gives attackers an increased possibility to generate more local events. However, it is assumed that event notifications are rate-limited in the implementation to address this threat.

9. MIB Considerations

SCTP-PF introduces new SCTP algorithms for failover and switchback with associated new state parameters. It is recommended that the SCTP-MIB defined in [RFC3873] is updated to support the management of the SCTP-PF implementation. This can be done by extending the sctpAssocRemAddrActive field of the SCTPAssocRemAddrTable to include information of the PF state of the destination address and by adding new fields to the SCTPAssocRemAddrTable supporting PotentiallyFailed.Max.Retrans (PFMR) and Primary.Switchover.Max.Retrans (PSMR) parameters.

10. IANA Considerations

This document does not create any new registries or modify the rules for any existing registries managed by IANA.

11. Acknowledgements

The authors wish to thank Michael Tuexen for his many invaluable comments and for his very substantial support with the making of this document.

12. Proposed Change of Status (to be Deleted before Publication)

Initially this work looked to entail some changes of the Congestion Control (CC) operation of SCTP and for this reason the work was proposed as Experimental. These intended changes of the CC operation have since been judged to be irrelevant and are no longer part of the specification. As the specification entails no other potential harmful features, consensus exists in the WG to bring the work forward as PS.

Initially concerns have been expressed about the possibility for the mechanism to introduce path bouncing with potential harmful network impacts. These concerns are believed to be unfounded. This issue is addressed in Appendix B.

It is noted that the feature specified by this document is implemented by multiple SCTP SW implementations and furthermore that various variants of the solution have been deployed in telephony signaling environments for several years with good results.

13. References

13.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.

13.2. Informative References

- [CARO02] Caro Jr., A., Iyengar, J., Amer, P., Heinz, G., and R. Stewart, "A Two-level Threshold Recovery Mechanism for SCTP", Tech report, CIS Dept, University of Delaware , 7 2002.
- [CARO04] Caro Jr., A., Amer, P., and R. Stewart, "End-to-End Failover Thresholds for Transport Layer Multi homing", MILCOM 2004 , 11 2004.
- [CARO05] Caro Jr., A., "End-to-End Fault Tolerance using Transport Layer Multi homing", Ph.D Thesis, University of Delaware , 1 2005.

- [FALLON08]
Fallon, S., Jacob, P., Qiao, Y., Murphy, L., Fallon, E.,
and A. Hanley, "SCTP Switchover Performance Issues in WLAN
Environments", IEEE CCNC 2008, 1 2008.
- [GRINNEMO04]
Grinnemo, K-J. and A. Brunstrom, "Performance of SCTP-
controlled failovers in M3UA-based SIGTRAN networks",
Advanced Simulation Technologies Conference , 4 2004.
- [IYENGAR06]
Iyengar, J., Amer, P., and R. Stewart, "Concurrent
Multipath Transfer using SCTP Multihoming over Independent
End-to-end Paths.", IEEE/ACM Trans on Networking 14(5), 10
2006.
- [JUNGMAIER02]
Jungmaier, A., Rathgeb, E., and M. Tuexen, "On the use of
SCTP in failover scenarios", World Multiconference on
Systemics, Cybernetics and Informatics , 7 2002.
- [NATARAJAN09]
Natarajan, P., Ekiz, N., Amer, P., and R. Stewart,
"Concurrent Multipath Transfer during Path Failure",
Computer Communications , 5 2009.
- [RFC3873] Pastor, J. and M. Belinchon, "Stream Control Transmission
Protocol (SCTP) Management Information Base (MIB)", RFC
3873, DOI 10.17487/RFC3873, September 2004,
<<http://www.rfc-editor.org/info/rfc3873>>.
- [RFC6458] Stewart, R., Tuexen, M., Poon, K., Lei, P., and V.
Yasevich, "Sockets API Extensions for the Stream Control
Transmission Protocol (SCTP)", RFC 6458, December 2011.

Appendix A. Discussions of Alternative Approaches

This section lists alternative approaches for the issues described in this document. Although these approaches do not require to update RFC4960, we do not recommend them from the reasons described below.

A.1. Reduce Path.Max.Retrans (PMR)

Smaller values for Path.Max.Retrans shorten the failover duration and in fact this is recommended in some research results [JUNGMAIER02] [GRINNEMO04] [FALLON08]. However to significantly reduce the failover time it is required to go down (as with PFMR) to Path.Max.Retrans=0 and with this setting SCTP switches to another

destination address already on a single timeout which may result in spurious failover. Spurious failover is a problem in [RFC4960] SCTP as the transmission of HEARTBEATS on the left primary path, unlike in SCTP-PF, is governed by 'HB.interval' also during the failover process. 'HB.interval' is usually set in the order of seconds (recommended value is 30 seconds) and when the primary path becomes inactive, the next HEARTBEAT may be transmitted only many seconds later. Indeed as recommended, only 30 secs later. Meanwhile, the primary path may since long have recovered, if it needed recovery at all (indeed the failover could be truly spurious). In such situations, post failover, an endpoint is forced to wait in the order of many seconds before the endpoint can resume transmission on the primary path and furthermore once it returns on the primary path the CWND needs to be rebuild anew - a process which the throughput already have had to suffer from on the alternate path. Using a smaller value for 'HB.interval' might help this situation, but it would result in a general waste of bandwidth as such more frequent HEARTBEATING would take place also when there are no observed troubles. The bandwidth overhead may be diminished by having the ULP use a smaller 'HB.interval' only on the path which at any given time is set to be the primary path, but this adds complication in the ULP.

In addition, smaller Path.Max.Retrans values also affect the 'Association.Max.Retrans' value. When the SCTP association's error count exceeds Association.Max.Retrans threshold, the SCTP sender considers the peer endpoint unreachable and terminates the association. Section 8.2 in [RFC4960] recommends that Association.Max.Retrans value should not be larger than the summation of the Path.Max.Retrans of each of the destination addresses. Else the SCTP sender considers its peer reachable even when all destinations are INACTIVE and to avoid this dormant state operation, [RFC4960] SCTP implementation SHOULD reduce Association.Max.Retrans accordingly whenever it reduces Path.Max.Retrans. However, smaller Association.Max.Retrans value decreases the fault tolerance of SCTP as it increases the chances of association termination during minor congestion events.

A.2. Adjust RTO related parameters

As several research results indicate, we can also shorten the duration of failover process by adjusting RTO related parameters [JUNGMAIER02] [FALLON08]. During failover process, RTO keeps being doubled. However, if we can choose smaller value for RTO.max, we can stop the exponential growth of RTO at some point. Also, choosing smaller values for RTO.initial or RTO.min can contribute to keep the RTO value small.

Similar to reducing Path.Max.Retrans, the advantage of this approach is that it requires no modification to the current specification, although it needs to ignore several recommendations described in the Section 15 of [RFC4960]. However, this approach requires to have enough knowledge about the network characteristics between end points. Otherwise, it can introduce adverse side-effects such as spurious timeouts.

The significant issue with this approach, however, is that even if the RTO.max is lowered to an optimal low value, then as long as the Path.Max.Retrans is kept at the [RFC4960] recommended value, the reduction of the RTO.max doesn't reduce the failover time sufficiently enough to prevent severe performance degradation during failover.

Appendix B. Discussions for Path Bouncing Effect

The methods described in the document can accelerate the failover process. Hence, they might introduce the path bouncing effect where the sender keeps changing the data transmission path frequently. This sounds harmful to the data transfer, however several research results indicate that there is no serious problem with SCTP in terms of path bouncing effect [CARO04] [CARO05].

There are two main reasons for this. First, SCTP is basically designed for multipath communication, which means SCTP maintains all path related parameters (CWND, ssthresh, RTT, error count, etc) per each destination address. These parameters cannot be affected by path bouncing. In addition, when SCTP migrates the data transfer to another path, it starts with the minimal or the initial CWND. Hence, there is little chance for packet reordering or duplicating.

Second, even if all communication paths between the end-nodes share the same bottleneck, the SCTP-PF results in a behavior already allowed by [RFC4960].

Appendix C. SCTP-PF for SCTP Single-homed Operation

For a single-homed SCTP association the only tangible effect of the activation of SCTP-PF operation is enhanced failure detection in terms of potential notification of the PF state of the sole destination address as well as, for idle associations, more rapid entering, and notification, of inactive state of the destination address and more rapid end-point failure detection. It is believed that neither of these effects are harmful, provided adequate dormant state operation is implemented, and furthermore that they may be particularly useful for applications that deploys multiple SCTP associations for load balancing purposes. The early notification of

the PF state may be used for preventive measures as the entering of the PF state can be used as a warning of potential congestion. Depending on the PMR value, the aggressive HEARTBEAT transmission in PF state may speed up the end-point failure detection (exceed of AMR threshold on the sole path error counter) on idle associations in case where relatively large HB.interval value compared to RTO (e.g. 30secs) is used.

Authors' Addresses

Yoshifumi Nishida
GE Global Research
2623 Camino Ramon
San Ramon, CA 94583
USA

Email: nishida@wide.ad.jp

Preethi Natarajan
Cisco Systems
510 McCarthy Blvd
Milpitas, CA 95035
USA

Email: prenatar@cisco.com

Armando Caro
BBN Technologies
10 Moulton St.
Cambridge, MA 02138
USA

Email: acar@bbn.com

Paul D. Amer
University of Delaware
Computer Science Department - 434 Smith Hall
Newark, DE 19716-2586
USA

Email: amer@udel.edu

Karen E. E. Nielsen
Ericsson
Kistavaegen 25
Stockholm 164 80
Sweden

Email: karen.nielsen@tieto.com

Network Working Group	M. Tuexen
Internet-Draft	I. Ruengeler
Updates: 4960 (if approved)	Muenster Univ. of Appl. Sciences
Intended status: Standards Track	R. Stewart
Expires: March 01, 2014	Adara Networks
	August 28, 2013

SACK-IMMEDIATELY Extension for the Stream Control Transmission Protocol
draft-ietf-tsvwg-sctp-sack-immediately-04.txt

Abstract

This document updates RFC 4960 by defining a method for the sender of a DATA chunk to indicate that the corresponding SACK chunk should be sent back immediately and not be delayed. It is done by specifying a bit in the DATA chunk header, called the I-bit, which can get set either by the SCTP implementation or by the application using an SCTP stack. Since unknown flags in chunk headers are ignored by SCTP implementations, this extension does not introduce any interoperability problems.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on March 01, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents

carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
2. Conventions	3
3. The I-bit in the DATA Chunk Header	3
4. Use Cases	4
4.1. Triggering at the Application Level	4
4.2. Triggering at the SCTP Level	4
5. Procedures	4
5.1. Sender Side Considerations	5
5.2. Receiver Side Considerations	5
6. Interoperability Considerations	5
7. Socket API Considerations	5
8. IANA Considerations	6
9. Security Considerations	7
10. Acknowledgments	7
11. References	7
11.1. Normative References	7
11.2. Informative References	7
Authors' Addresses	7

1. Introduction

According to [RFC4960] the receiver of a DATA chunk should use delayed SACKs. This delaying is completely controlled by the receiver of the DATA chunk and remains the default behavior.

In specific situations the delaying of SACKs results in reduced performance of the protocol:

1. If such a situation can be detected by the receiver, the corresponding SACK can be sent immediately. For example, [RFC4960] recommends the immediate sending if the receiver has detected message loss or message duplication.

2. However, if the situation can only be detected by the sender of the DATA chunk, [RFC4960] provides no method of avoiding a delay in sending the SACK. Examples of these situations include ones which require interaction with the application (e.g. applications using the SCTP_SENDER_DRY_EVENT, see Section 4.1) and ones which can be detected by the SCTP stack itself (e.g. closing the association, hitting window limits or resetting streams, see Section 4.2).

To overcome the limitation described in the second case, this document describes a simple extension of the SCTP DATA chunk by defining a new flag, the I-bit. The sender of a DATA chunk indicates by setting this bit that the corresponding SACK chunk should not be delayed.

2. 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. The I-bit in the DATA Chunk Header

The following Figure 1 shows the extended DATA chunk.

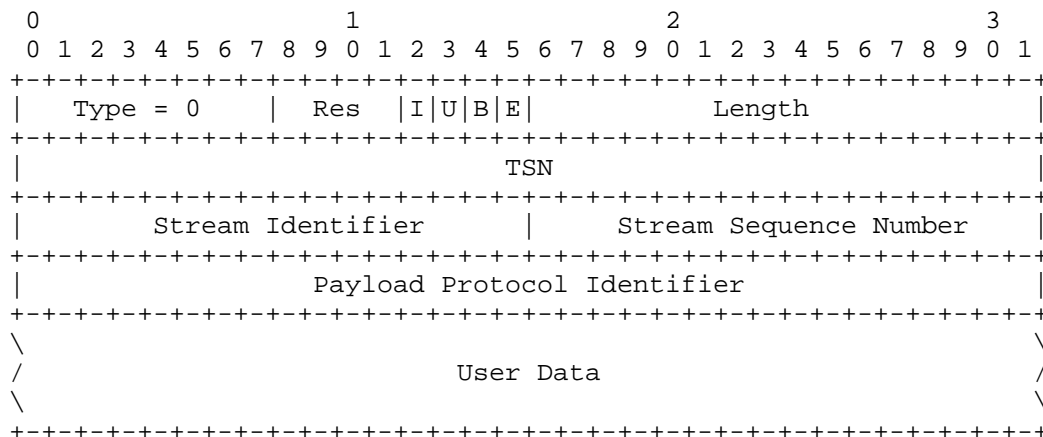


Figure 1: Extended DATA chunk format

The only difference between the DATA chunk in Figure 1 and the DATA chunk defined in [RFC4960] is the addition of the I-bit in the flags field of the DATA chunk header.

This bit was Reserved in [RFC4960]. [RFC4960] specified that this bit should be set to 0 by the sender and ignored by the receiver.

4. Use Cases

The setting of the I-bit can either be triggered by the application using SCTP or by the SCTP stack itself. The following two subsections provide a non-exhaustive list of examples.

4.1. Triggering at the Application Level

One example of a situation in which it may be desirable for an application to trigger setting of the I-bit involves the SCTP_SENDER_DRY_EVENT in the SCTP socket API [RFC6458]. Upper layers of SCTP using the socket API as defined in [RFC6458] may subscribe to the SCTP_SENDER_DRY_EVENT for getting a notification as soon as no user data is outstanding anymore. To avoid an unnecessary delay while waiting for such an event, the application can request the setting of the I-Bit when sending the last user message before waiting for the event. This results in setting the I-bit of the last DATA chunk corresponding to the user message and is possible using the extension of the socket API described in Section 7.

4.2. Triggering at the SCTP Level

There are also situations in which the SCTP implementation can set the I-bit without interacting with the upper layer.

If the association is in the SHUTDOWN-PENDING state, setting the I-bit reduces the number of simultaneous associations for a busy server handling short living associations.

Another case is where the sending of a DATA chunk fills the congestion or receiver window. Setting the I-bit in these cases improves the throughput of the transfer.

If an SCTP association supports the SCTP Stream Reconfiguration extension defined in [RFC6525], the performance can be improved by setting the I-bit when there are pending reconfiguration requests that require that there be no outstanding DATA chunks.

5. Procedures

5.1. Sender Side Considerations

Whenever the sender of a DATA chunk can benefit from the corresponding SACK chunk being sent back without delay, the sender MAY set the I-bit in the DATA chunk header. Please note that it is irrelevant to the receiver why the sender has set the I-bit.

Reasons for setting the I-bit include, but are not limited to, the following (see Section 4 for the benefits):

- o The application requests to set the I-bit of the last DATA chunk of a user message when providing the user message to the SCTP implementation (see Section 7).
- o The sender is in the SHUTDOWN-PENDING state.
- o The sending of a DATA chunk fills the congestion or receiver window.
- o The sending of an Outgoing SSN Reset Request Parameter or an SSN/TSN Reset Request Parameter is pending, if the association supports the Stream Reconfiguration extension defined in [RFC6525].

5.2. Receiver Side Considerations

On reception of an SCTP packet containing a DATA chunk with the I-bit set, the receiver SHOULD NOT delay the sending of the corresponding SACK chunk, i.e., the receiver SHOULD immediately respond with the corresponding SACK chunk.

6. Interoperability Considerations

According to [RFC4960] the receiver of a DATA chunk with the I-bit set should ignore this bit when it does not support the extension described in this document. Since the sender of the DATA chunk is able to handle this case, there is no requirement for negotiating the support of the feature described in this document.

7. Socket API Considerations

This section describes how the socket API defined in [RFC6458] is extended to provide a way for the application to set the I-bit.

Please note that this section is informational only.

A socket API implementation based on [RFC6458] needs to be extended to allow the application to set the I-bit of the last DATA chunk when sending each user message.

This can be done by setting a flag called `SCTP_SACK_IMMEDIATELY` in the `snd_flags` field of the struct `sctp_sndinfo` structure when using `sctp_sendv()` or `sendmsg()`. If the deprecated struct `sctp_sndrcvinfo` structure is used instead when calling `sctp_send()`, `sctp_sendx()`, or `sendmsg()`, the `SCTP_SACK_IMMEDIATELY` flag can be set in the `sinfo_flags` field. When using the deprecated function `sctp_sendmsg()` the `SCTP_SACK_IMMEDIATELY` flag can be in the `flags` parameter.

8. IANA Considerations

[NOTE to RFC-Editor:

"RFCXXXX" is to be replaced by the RFC number you assign this document.

]

Following the chunk flag registration procedure defined in [RFC6096], IANA should register a new bit, the I-bit, for the DATA chunk. The suggested value is 0x08 and the reference should be RFCXXXX.

This requires an update of the "DATA Chunk Flags" registry for SCTP:

DATA Chunk Flags

Chunk Flag Value	Chunk Flag Name	Reference
0x01	E bit	[RFC4960]
0x02	B bit	[RFC4960]
0x04	U bit	[RFC4960]
0x08	I Bit	[RFCXXXX]
0x10	Unassigned	
0x20	Unassigned	
0x40	Unassigned	
0x80	Unassigned	

9. Security Considerations

See [RFC4960] for general security considerations for SCTP. In addition, a malicious sender can force its peer to send packets containing a SACK chunk for each received packet containing DATA chunks instead of every other. This could impact the network, resulting in more packets sent on the network, or the peer because the generating and sending of the packets has some processing cost. However, the additional packets can only contain the most simplest SACK chunk (no gap reports, no duplicate TSNs), since in case of packet drop or reordering in the network a SACK chunk would be sent immediately anyway. Therefore this does neither introduce a significant additional processing cost on the receiver side. This does not result in more traffic in the network than a receiver that sends a SACK for every packet, which is already permitted.

10. Acknowledgments

The authors wish to thank Mark Allmann, Brian Bidulock, David Black, Anna Brunstrom, Gorry Fairhurst, Janardhan Iyengar, Kacheong Poon, and Michael Welzl for their invaluable comments.

11. References

11.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.
- [RFC6096] Tuexen, M. and R. Stewart, "Stream Control Transmission Protocol (SCTP) Chunk Flags Registration", RFC 6096, January 2011.

11.2. Informative References

- [RFC6458] Stewart, R., Tuexen, M., Poon, K., Lei, P., and V. Yasevich, "Sockets API Extensions for the Stream Control Transmission Protocol (SCTP)", RFC 6458, December 2011.
- [RFC6525] Stewart, R., Tuexen, M., and P. Lei, "Stream Control Transmission Protocol (SCTP) Stream Reconfiguration", RFC 6525, February 2012.

Authors' Addresses

Michael Tuexen
Muenster University of Applied Sciences
Stegerwaldstr. 39
48565 Steinfurt
DE

Email: tuexen@fh-muenster.de

Irene Ruengeler
Muenster University of Applied Sciences
Stegerwaldstr. 39
48565 Steinfurt
DE

Email: i.ruengeler@fh-muenster.de

Randall R. Stewart
Adara Networks
Chapin, SC 29036
US

Email: randall@lakerest.net

Network Working Group
Internet-Draft
Intended status: Informational
Expires: January 7, 2014

C. Lai
W. Wang
S. Yang
T. Eckert
F. Yip
Cisco Systems
July 6, 2013

Normalization Marker for AF PHB Group in DiffServ
draft-lai-tsvwg-normalizer-02

Abstract

In DiffServ, preferential dropping of packets in AF PHB groups has long been considered beneficial, typically for video flows with discardable packets. Unfortunately, the ecosystem of bandwidth contention at congestion is very likely to discourage those video endpoints from generating packets with lower precedence markings, i.e. they would lose more packets if doing so. Thus, to offer an incentive for more collaborative and mutually beneficial behaviors of video endpoints in AF PHB groups, we propose a Normalization Marker (NM) for traffic conditioning at network edges. Deployment of NM will encourage the video endpoints to generate finer layers of intra-flow precedence (IFP) with discardable packets in more balanced distributions.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on January 7, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	3
2. Background	6
2.1. Video Packets in Structure	6
2.2. Intra-Flow Precedence (IFP)	8
2.3. Mapping IFP to AF Markings	9
3. Normalization Marker (NM)	11
3.1. Color-Aware vs. Color-Blind Mode	12
3.2. Distribution Meter	12
3.3. Normalizer	13
4. Acknowledgements	13
5. IANA Considerations	13
6. Security Considerations	13
7. References	13
7.1. Normative References	13
7.2. Informative References	14
Authors' Addresses	14

1. Introduction

Assured Forwarding (AF) Per-Hop Behavior (PHB) groups are described in [RFC2597] (with terminology clarified in [RFC3260]) for DiffServ (DS) multimedia service classes such as realtime video conferencing and on-demand streaming. Four AF PHB groups have been defined in [RFC4594] with DS codepoint (DSCP): AF1x, AF2x, AF3x and AF4x where x=1, 2 or 3 for drop precedence in each independent AF PHB group. The DS nodes that support an AF PHB group must set configuration of Active Queue Management (AQM) properly w.r.t. those DSCP markings. For example, for AF4x PHB group which includes AF41, AF42 and AF43 markings, an AQM implementation by Weighted Random Early Detection (WRED) should be configured with some drop probabilities and queue thresholds such that the packet loss rate of AF41 \leq AF42 \leq AF43 on congestion of the queue.

For an AF PHB group, a DS boundary node or host in the DS domain should use a marking algorithm that properly assigns AF markings of drop precedence to all packets w.r.t. the traffic profiles and Service Level Agreements (SLA). For example, [RFC2697] and [RFC2698] use a token-bucket mechanism for metering each stream of packets and respectively define "srTCM" and "trTCM" markers, to mark packets by the data rate and burst size limit in traffic profiles. Those rate-control markers can be useful at DS boundary nodes for traffic conditioning [RFC2475] and to support IntServ/RSVP traffic over DS regions [RFC2998]. Multiple markers may be applied to the same stream, either on the same or multiple DS nodes along the path. For example, srTCM and trTCM can operate in a so-called "color-aware" mode such that for each incoming packet that already carries an AF marking, the local srTCM/trTCM either keeps the same or lowers the drop precedence of that packet by metering.

However, modern video codec technologies are being advanced not only in coding efficiency (i.e. better compression ratio) but also in two key areas for transport on IP networks: (1) encoder rate-control and dynamic adaptation; (2) ability to generate discardable packets in multiple layers to tolerate packet losses in the network without significant degradation of video quality observed at the decoder. For (1), the encoder dynamically limits its output rate of packets into the AF PHB group, i.e., the encoder's host is the first DS node equipped with srTCM/trTCM if it marks packets in that behavior. The next DS node is the first-hop router which may add extra srTCM/trTCM to enforce the traffic conditioning or policing from the network's perspective. Thus, we consider this an incentive for (1) because an encoder using a self rate-control is less likely to see packet losses by the network. Unfortunately, an incentive for (2) is arguably missing today.

To see the missing incentive for (2), consider the following example where 2 video flows A and B with rate control are sent in AF4x PHB group. Each sends 5Mbps on average with some burstiness, but still complies with the rate and burst limit in its traffic profile. However, A and B generate packets with AF4x markings in different distributions of percentage:

Flow A

80% or 4Mbps in AF41

20% or 1Mbps in AF42

0% or 0Mbps in AF43

Flow B

40% or 2Mbps in AF41

40% or 2Mbps in AF42

20% or 1Mbps in AF43

Flow B at above is likely using a more advanced video technology to generate multiple layers of discardable video packets, and thus, its distribution of AF4x markings looks finer and more balanced. That is, flow B acts more friendly to other flows in this AF4x PHB group.

Thus, we argue that the ecosystem in practical deployment should offer an incentive for flows to behave similarly to what flow B is doing above, i.e., on congestion, the AF4x PHB group should try to drop packets in the same amount from each flow, while a flow with finer layers of discardable packets and/or in a more balanced distribution should be able to benefit from its own efforts and see good results in video quality preservation.

Unfortunately, this incentive is still missing today. Suppose that congestion occurs in the AF4x WRED queue where A and B compete for bandwidth and there is no other flow, for simplicity. B's packet loss rate is very likely to become higher than A's, despite B's effort of acting friendly:

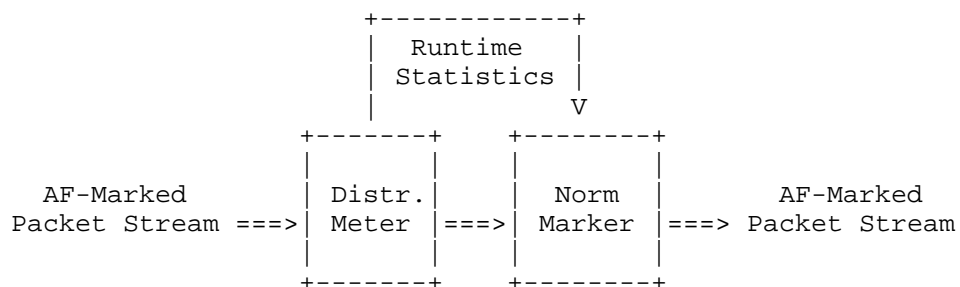
- o If the queue drops 1Mbps in total,

A sees 0% or 0Mbps loss;

B sees 20% or 1Mbps loss (all its AF43 are lost).

- o If the queue drops 4Mbps in total,
 - A sees 20% or 1Mbps loss (all its AF42 are lost);
 - B sees 60% or 3Mbps loss (all its AF42 and AF43 are lost).

Thus, to create the missing incentive at above, we propose a new "Normalization Marker" (NM) and describe it in this memo. NM can be deployed on DS boundary nodes for traffic conditioning in practical deployment with AF PHB groups for multimedia service classes. In summary, if NM is applied to a DS boundary node for an AF PHB group, it re-assigns the AF markings of all packets per flow such that the distributions of the AF markings are similar in all flows, i.e., it "normalizes" the distributions of AF markings in all flows. It also attempts to maintain the original orders of the intra-flow drop precedence carried by the input AF markings, as linearly as possible. After the AF-marking distributions are normalized, all those flows should see very similar packet loss rates at AQM for this AF PHB group on congestion of the queue. Then, a codec implementation may have better video quality preservation on network congestion if it employs a more advanced video technology to generate discardable packets with finer markings of drop precedence in a more balanced distribution.



Normalization Marker (NM) with AF PHB Group

Figure 1

Note that the use of NM is not necessarily limited to video service classes, but could be extended to wherever AF PHB groups can be used, or to any other PHB groups that require a similar incentive NM can provide.

2. Background

2.1. Video Packets in Structure

Modern video codec technologies such as ITU-T H.264/MPEG-4 AVC [H264] typically generate a stream of encoded video packets with internal structure of data dependency for decoding. This has been designed for at least 3 fundamental reasons:

- o **Coding Efficiency:** An encoder improves its coding efficiency typically by reducing spatial and temporal redundancy of the input. For video, spatial redundancy is reduced by intra-frame motion prediction and compensation, while temporal redundancy refers to inter-frame since a video stream is composed of a sequence of frames or pictures in the temporal order. With motion prediction, a frame can be encoded by referencing some pixels of the picture data that will be decoded earlier either in the same (intra) or another (inter) frame so that it can use significantly fewer bits to encode this frame. The frame where the pixels are referenced by any other frame is thus called a referenced frame in the video stream; for example, Instantaneous Decoding Refresh (IDR) in H.264 or Intra (I) frames are typically referenced by subsequential frames, while Predictive (P) frames may be referenced at the encoder's choice, by the Group of Picture (GOP) profile, and/or by some proprietary algorithm in the codec implementation.
- o **Lossy Network:** To use network transport that may lose packets, the encoder may choose to generate a stream with two or more layers each of which the packets are marked with some layer identifier (ID). The network can simply use the layer ID to determine the drop precedence of each packet in the video stream.
 - * **Layers in Hierarchy of Dependency:** If these layers are coded in hierarchy of dependency, the packets in an "enhancement" layer will depend on 1 or more "base" layers to get decoded without errors, while packets in a base layer without dependency can be independently decoded without errors.
 - + If some enhancement layer packets are lost, the decoding errors in that picture frame will not stay or cascade to other frames given that no others depend on those lost data. This nice property allows the network to safely drop packets in some enhancement layers, if needed, without badly impacting the video quality at decoder.
 - + If some base layer packets are lost, the impact can be severe since these decoding errors will stay in buffer and

cascade to all other picture pixels that depend on the lost data to decode in the current and/or a later frame. This impact can last tens of seconds as the video quality continues getting worse, resulting in unpleasant user experiences, until the decoder receives the next IDR or I frame, either on-demand or periodically, to remove those errors.

For example, H.264 Annex G defines Scalable Video Coding (SVC) using a 3-dimensional (i.e. spatical, temporal and quality) hierarchy of layer dependency at the encoder's choice, but for simplicity, it also defines a scalar number called Priority ID (PID) in its header so the network could instead use PID, if set by the encoder, to determine drop precedence in the stream.

- * Layers NOT in Hierarchy of Dependency: Sometimes the encoder will generate multiple layers without any dependency between those layers. These mechanisms usually enlarge the amount of encoded video data for vairous purposes. For example,
 - + Forward Error Correction (FEC) may be used at the encoder to generate extra FEC packets, so that the decoder can tolerate certain amounts of packet losses.
 - + Simulcast (i.e. simultaneous multicast) by an encoder will actually generate multiple layers each of which can be transmitted and decoded independently, in parallel by IP or application multicast. Each layer carries video in a different resolution and/or quality. The decoder can choose 1 or more of those layers to receive according to the required, available or detected bandwidth, packet losses, delays, jitter etc. in its network service.

With FEC and/or Simulcast, the encoder can still mark the packets with different drop precedence in those layers to better protect the more important data for video quality at decoding when congestion occurs.

- o In-Band Signaling: An encoded video stream usually carries in-band control messages that are most critical for adequate encoder and decoder behaviors. For example,
 - * H.264 Annex D defines Supplemental Enhancement Information (SEI), which could also carry proprietary codec parameters. These in-band control signals should be given the highest drop precedence.

- * Real Time Control Protocol (RTCP) carries in-band control messages for Real Time Protocol (RTP) [RFC3550], which is mostly used for realtime multimedia transmission on IP networks. RTCP messages are defined as RTP packets with special payload types in the RTP stream. RTCP packets should be given the highest drop precedence but should receive the same delay/jitter as regular RTP packets in the same stream.

2.2. Intra-Flow Precedence (IFP)

For abstraction, we define "Intra-Flow Precedence" (IFP) to represent the drop precedence in one individual flow that may carry a video stream of IP packets in multimedia networks. Here is a summary of IFP characteristics:

- o IFPs are drop precedence levels that are only significant within each individual flow.
- o IFPs are integer numbers that can be numerically compared if needed. 0 represents the highest precedence. The larger numerical value an IFP is, the lower precedence it represents.
- o The number of IFP levels in each flow is not necessarily the same.
- o IFPs between any 2 flows should NOT be compared to determine drop precedence between their packets in a queue.
- o IFPs may be assigned by the original encoder of the stream and carried in some bits field of all packets in the stream.
- o IFPs may be assigned or re-assigned by a middle box or router if it is capable of understanding the stream packet format and codec semantics.

For example, an H.264 AVC flow may have the following IFP assignments at the video encoder's choice.

IFP = 0 for in-band signals

IFP = 1 for IDR frames

IFP = 2 for referenced P (rP) frames

IFP = 3 for non-referenced P (nrP) frames and others

IFP assignments as well as their distribution can vary a lot among different encoder implementations and codec profiles. For example, some encoders may generate both long-term and short-term referenced P

frames, where a long-term referenced P frame should have higher drop precedence. In case of H.264 SVC, the IFP assignments could simply be the same as the PID assignments if set by the encoder properly, or be calculated based on the SVC layer ID that has 3 tuples for the spatial, temporal and quality dimensions, respectively.

2.3. Mapping IFP to AF Markings

When a flow is sent in an AF PHB group, the number of its IFP levels is not necessarily equal to the number of the AF markings. In fact, since each of the currently defined AF PHB groups has only 3 AF markings, it is likely that an encoder or DS node needs to apply an n-to-1 mapping from IFPs to AF markings in practice.

The mapping decision is made usually by the encoder, but can also be made by another DS node if necessary and if the DS node is able to understand the encoded video packets, which may require Deep Packet Inspection (DPI), e.g. to read in RTP payload and parse the H.264 headers [RFC6184], or in a proprietary bits field in the IP payload, to retrieve or calculate the IFP of each packet in a flow before locally mapping the IFP to an AF marking.

This n-to-1 mapping can be arbitrary but should be appropriate. Consider 2 IFPs, say x and y , where x and y are mapped to AF markings $AF(x)$ and $AF(y)$, respectively. Then, the mapping should ideally obey the following criteria to keep linearity from IFPs to AF markings.

If $x < y$, $AF(x) \leq AF(y)$;

If $x > y$, $AF(x) \geq AF(y)$.

Although the above two do NOT imply that if $x = y$, $AF(x) = AF(y)$, it is usually so in practical implementation as it is straightforward. Then, if the encoder algorithm generates a lot of packets with the same IFP, all those packets will be assigned the same AF marking, possibly resulting in an unbalanced distribution of AF markings in the AF PHB group. Thus, an encoder with advanced technologies should make good efforts to generate packets with a finer and more balanced IFP distribution in the first place.

For example, if AF4x PHB group is used to send an H.264 AVC flow with the IFP assignments in the example of Section 2.2, one possible IFP-to-AF4x mapping is:

$AF(0) = AF41$

$AF(1) = AF41$

AF(2) = AF42

AF(3) = AF43

This mapping actually results in the following AF markings:

AF41 for in-band signals and IDR frames

AF42 for referenced P (rP) frames

AF43 for non-referenced P (nrP) frames and others

Now, consider two encoders that generate flow A and B, respectively, both using this mapping, but with different IFP distributions as follows.

Flow A

5% in IFP=0 for in-band signals

75% in IFP=1 for IDR frames

20% in IFP=2 for rP frames

Flow B

5% in IFP=0 for in-band signals

35% in IFP=1 for IDR frames

40% in IFP=2 for rP frames

20% in IFP=3 for nrP frames

Thus,

Flow A

80% in AF41

20% in AF42

0% in AF43

Flow B

40% in AF41

40% in AF42

20% in AF43

This results in exactly the two AF marking distributions that we have previously used in Section 1.

Note that in terms of encoded data size, an IDR frame is typically 10 times larger than a P frame on average. Assume that flow B's coding efficiency has rP twice as large as nrP. Then, flow A and B might be sending frames periodically in patterns by Group of Picture (GOP) as follows:

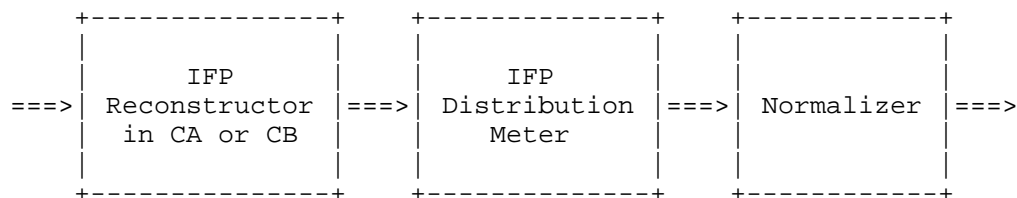
Flow A: IDR, rP, rP, rP

Flow B: IDR, rP, nrP, rP, nrP, rP, nrP, rP, nrP

If so, it shows that flow B's encoder is making efforts to generate discardable packets with more layers in a more balanced distribution, which is desirable.

3. Normalization Marker (NM)

Referring to Figure 2, NM has 3 major components: IFP reconstructor, IFP distribution meter, and normalizer. NM may operate in either "color-aware" (CA) or "color-blind" (CB) mode.



Normalization Marker (NM) Architecture

Figure 2

The packets arrive at the IFP reconstructor which determines the IFP of each packet depending on whether NM is in CA or CB mode. This is fed into the IFP distribution meter that keeps a runtime statistics. Then, by the runtime statistics and the IFP of the very packet, the normalizer writes a proper AF-marking in that packet.

3.1. Color-Aware vs. Color-Blind Mode

When NM operates in "color-aware" (CA) mode, it reads the incoming AF-markings that are carried in the packets as the drop precedence. This CA mode should be supported in all NM implementations.

When NM operates in "color-blind" (CB) mode, which is optionally supported, it reads certain bits field(s) other than the AF-markings in the packets to determine the actual drop precedence of that packet. This implies that NM may need DPI in the packets, e.g. parsing into H.264 AVC header in each RTP packets, or alternatively use some method where the drop precedence is carried from the encoder in a customized bits field other than the AF-marking in each packet.

In comparison, CB is more complex than CA in implementation. However, CB could probably produce better normalization results because the AF-markings are actually outcomes of an n-to-1 mapping from IFPs, as previously mentioned in Section 2.3, which can reduce granularity, e.g. for IFPs x and y , if $x > y$ at encoder, it is possible that $AF(x) = AF(y)$ when NM sees those packets in CA mode. On the contrary, NM in CB mode may reconstruct IFPs $x > y$ for those packets by local DPI.

Note that NM in CB mode may fail to determine the IFP of a packet for various reasons at runtime. If so, NM should randomly assign an IFP to each of those packets with an even distribution over the IFPs. The failure could be due to payload encryption that prevents DPI. Another reason may be that the NM does not support the codec used for encoding those packets in the flow. For example, an NM might only support H.264 AVC but is unable to parse packets in H.264 Annex G (SVC), so it fails to determine the IFPs of packets in an H.264 SVC flow.

3.2. Distribution Meter

The IFP distribution meter keeps a runtime statistics of the IFPs per flow so that the normalizer will be able to assign a proper AF-marking for each packet. The types of statistics to collect at runtime depend on the NM algorithm in the implementation.

For example, an NM implementation may keep a counter of packets per IFP in a flow since the beginning of the flow's lifetime. Another implementation may choose to keep only the running average of the packet counter per IFP. An even simpler implementation may choose to keep only the running average of IFPs of all packets per flow.

3.3. Normalizer

The normalizer should reference the runtime statistics kept by the IFP distribution meter, and adaptively map the IFP of the very packet to an AF marking, such that the resulting AF-marking distributions for all flows are similar or even identical to a target distribution.

The target distribution of an NM can be simply an even distribution over all possible AF-markings in the AF PHB group. However, in a more complex NM implementation, it may allow configuration for other target distributions as appropriate with the AQM configuration.

4. Acknowledgements

The authors would like to thank many colleagues for comments and supports, and thank Shuai Dai for testing NM with actual H.264 video endpoints.

5. IANA Considerations

This memo includes no request to IANA.

6. Security Considerations

This memo has no security consideration at the time of writing.

7. References

7.1. Normative References

- [RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and W. Weiss, "An Architecture for Differentiated Services", RFC 2475, December 1998.
- [RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", RFC 2597, June 1999.
- [RFC2697] Heinanen, J. and R. Guerin, "A Single Rate Three Color Marker", RFC 2697, September 1999.
- [RFC2698] Heinanen, J. and R. Guerin, "A Two Rate Three Color Marker", RFC 2698, September 1999.
- [RFC2998] Bernet, Y., Ford, P., Yavatkar, R., Baker, F., Zhang, L.,

Speer, M., Braden, R., Davie, B., Wroclawski, J., and E. Felstaine, "A Framework for Integrated Services Operation over Diffserv Networks", RFC 2998, November 2000.

[RFC3260] Grossman, D., "New Terminology and Clarifications for Diffserv", RFC 3260, April 2002.

[RFC4594] Babiarz, J., Chan, K., and F. Baker, "Configuration Guidelines for DiffServ Service Classes", RFC 4594, August 2006.

7.2. Informative References

[H264] ITU-T Recommendation H.264, "Advanced video coding for generic audiovisual services", March 2010.

[RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.

[RFC6184] Wang, Y., Even, R., Kristensen, T., and R. Jesup, "RTP Payload Format for H.264 Video", RFC 6184, May 2011.

Authors' Addresses

Cheng-Jia Lai
Cisco Systems
170 West Tasman Drive
San Jose, CA 95134
US

Email: chelai@cisco.com

Wenyi Wang
Cisco Systems
170 West Tasman Drive
San Jose, CA 95134
US

Email: wenywang@cisco.com

Stan Yang
Cisco Systems
170 West Tasman Drive
San Jose, CA 95134
US

Email: stanyang@cisco.com

Toerless Eckert
Cisco Systems
170 West Tasman Drive
San Jose, CA 95134
US

Email: eckert@cisco.com

Fred Yip
Cisco Systems
San Diego, CA
US

Email: fyip@cisco.com

Transport Area
Internet-Draft
Intended status: Informational
Expires: June 7, 2014

T. Moncaster, Ed.
J. Crowcroft
University of Cambridge
M. Welzl
University of Oslo
D. Ros
Telecom Bretagne
M. Tuexen
Muenster Univ. of Appl. Sciences
December 4, 2013

Problem Statement: Why the IETF Needs Defined Transport Services
draft-moncaster-tsvwg-transport-services-01

Abstract

The IETF has defined a wide range of transport protocols over the past three decades. However, the majority of these have failed to find traction within the Internet. This has left developers with little choice but to use TCP and UDP for most applications. In many cases the developer isn't interested in which transport protocol they should use. Rather they are interested in the set of services that the protocol provides to their application. TCP provides a very rich set of transport services, but offers no flexibility over which services can be used. By contrast, UDP provides a minimal set of services.

As a consequence many developers have begun to write application-level transport protocols that operate on top of UDP and offer them some of the flexibility they are looking for. We believe that this highlights a real problem: applications would like to be able to specify the services they receive from the transport protocol, but currently transport protocols are not defined in this fashion. There is an additional problem relating to how to ensure new protocols are able to be adopted within the Internet, but that is beyond the scope of this problem statement.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on June 7, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	2
1.1. Changes in This Version (to be removed by RFC Editor) . .	3
2. Transport Services	3
2.1. Identifying Transport Services	4
2.2. Exposing Transport Services	4
3. Why Now?	5
4. Security Considerations	6
5. IANA Considerations	6
6. Conclusions	6
7. Contributors and Acknowledgements	7
8. Comments Solicited	7
9. References	7
9.1. Normative References	7
9.2. Informative References	8

1. Introduction

The IETF has defined a wide array of transport protocols including UDP [RFC0768], TCP [RFC0793], SCTP [RFC4960], UDP-Lite [RFC3828], DCCP [RFC4340] and MPTCP [RFC6824]. In most cases new protocols have been defined because the IETF has established that there is a need for a set of behaviours than cannot be offered by any existing transport protocol.

However, for an application programmer, using protocols other than TCP or UDP can be hard: not all protocols are available everywhere, hence a fall-back solution to TCP or UDP must be implemented. Some protocols provide the same services in different ways. Layering decisions must be made (e.g. should a protocol be used natively or over UDP?). Because of these complications, programmers often resort to either using TCP (even if there is a mismatch between the services provided by TCP and the services needed by the application) or implementing their own customised solution over UDP, and the opportunity of benefiting from other transport protocols is lost. Since all these protocols were developed to provide services that solve particular problems, the inability of applications to make use of them is in itself a problem. Implementing a new solution e.g. over UDP also means re-inventing the wheel (or, rather, re-implementing the code) for a number of general network functions such as methods to interoperate through NATs and PMTUD.

We believe this mismatch between the application layer and transport layer can be addressed in a simple fashion. If an API allowed applications to request transport services without specifying the protocol, the transport system underneath could automatically try to make the best of its available resources. It could use available transport protocols in a way that is most beneficial for applications and without the application needing to worry about problems with middlebox traversal. Adopting this approach could give more freedom for diversification to designers of Operating Systems.

1.1. Changes in This Version (to be removed by RFC Editor)

From draft-moncaster-tsvwg-transport-services-00 to -01: Editorial corrections and clarifications including:

- * Updated Section 2.1 to highlight that we will take a hybrid approach to identifying Transport Services, both top down (by examining existing APIs) and bottom up (by looking at existing transport protocols).
- * Updated Section 2.2 to commit to delivering at least one example API for this work.
- * Replaced Section 4. The new version makes it clear that we will preserve the status quo where the transport may or may not choose to implement security.

2. Transport Services

The transport layer provides many services both to the end application (e.g. multiplexing, flow control, ordering, reliability)

and to the network (e.g. congestion control). For the purposes of this document we define Transport Services as follows:

- o A Transport Service is any service provided by the transport layer that can only be correctly implemented with information from the application.

The key word here is "information" -- many existing transport protocols function perfectly adequately because the choice of protocol implicitly includes information about the desired transport capabilities. For instance the choice of TCP implies a desire for reliable, in-order data delivery. However we think that such implicit information is not always sufficient. The rest of this section explains how we propose to identify Transport Services and how those services might then be exposed to the application.

2.1. Identifying Transport Services

One of the key aspects of this work is how to identify which Transport Services should actually be supported. We are taking a two-pronged approach. Rather than trying to identify every possible service that popular applications might need, we will survey a given set of common APIs that applications use to communicate across the network. We will complement this with a bottom-up approach where we establish the set of services that have already been published in RFCs coming from the Transport Area. This way, much of the discussion about the need to specify these services has already taken place, and it is unnecessary to re-visit those discussions. It is our hope that this approach will lead to identifying a set of service primitives that can be combined to offer a rich set of services to the application.

2.2. Exposing Transport Services

These Transport Services would be exposed to the application via an API. The definition of such an API and the functionality underneath the API are beyond the scope of this problem statement. We briefly describe three possible approaches below.

One approach could be to develop a transport system that fully operates inside the Operating System. This transport system would provide all the defined services for which it can use TCP as a fall-back at the expense of efficiency (e.g., TCP's reliable in-order delivery is a special case of reliable unordered delivery, but it may be less efficient). To test whether a particular transport is available it could take the Happy Eyeballs [I-D.wing-tsvwg-happy-eyeballs-sctp] approach proposed for SCTP -- if the SCTP response arrives too late then the connection just uses TCP

and the SCTP association information could be cached so that a future connection request to the same destination IP address can automatically use it.

Polyversal TCP [PVTCP] offers another possible approach. This starts by opening a TCP connection and then attempts to establish other paths using different transports. The TCP connection ensures there's always a stable fallback. Having established the initial connection, PVTCP can then use service requests coming through `setsockopt()` to select the most appropriate transport from the available set.

Another approach could be to always rely on UDP only, and develop a whole new transport protocol above UDP which provides all the services, using a single UDP port. Instead of falling back to TCP, this transport system could return an error in case there is no other instance of the transport system available on the other side; the first packets could be used to signal which service is being requested to the other side (e.g., unordered delivery requires the receiving end to be aware of it).

3. Why Now?

So why do we need to deal with this issue now? There are several answers. Firstly, after several decades of dominance by various flavours of TCP and UDP (plus limited deployment of SCTP [RFC4960]), transport protocols are undergoing significant changes. Recent standards allow for parallel usage of multiple paths (MPTCP [RFC6824] and CMT-SCTP [I-D.tuexen-tsvwg-sctp-multipath]) while other standards allow for scavenger-type traffic (LEDBAT [RFC6817]). What sets these apart from e.g. DCCP [RFC4340] is that they have already seen deployment in the wild -- one of the Internet's most popular applications, BitTorrent, uses LEDBAT and MPTCP is already seeing deployment in major operating systems [Bonaventure-Blog]. Meanwhile there is a trend towards tunnelling transports inside UDP -- SCTP over DTLS over UDP is now being shipped with a popular browser in order to support WebRTC [RFC6951][I-D.ietf-tsvwg-sctp-dtls-encaps] while RTMFP [I-D.thornburgh-adobe-rtmfp] and QUIC [QUIC] are recent examples of transport protocols that are implemented over UDP in user space. In a similar vane, Minion [I-D.iyengar-minion-protocol] is a proposal to realise some SCTP-like services with a downwards-compatible extension to TCP.

All of a sudden, application developers are faced with a heterogeneous, complex set of protocols to choose from. Every protocol has its pro's and con's, but often the reasons for making a particular choice depend not on the application's preferences but on the environment (e.g., the choice of Minion vs. SCTP would depend on whether SCTP could successfully be used on a given network path).

Choosing a protocol that isn't guaranteed to work requires implementing a fall-back method to e.g. TCP, and making the best possible choice at all times may require sophisticated network measurement techniques. The process could be improved by using a cache to learn which protocols previously worked on a path, but this wouldn't always work in a cloud environment where virtual machines can and do migrate between physical nodes.

We therefore argue that it is necessary to provide mechanisms that automate the choice and usage of the transport protocol underneath the API that is exposed to applications. As a first step towards such automation, we need to define the services that the transport layer should expose to an application (as opposed to today's typical choice of TCP and UDP).

4. Security Considerations

Whether or not to enable TLS[RFC5246] is currently left up to individual protocol implementations to decide. While there is some debate about whether this is correct we have chosen to keep the status quo.

5. IANA Considerations

This document makes no request to IANA although in future an IANA register of Transport Services may be required.

6. Conclusions

After decades of relative stagnation the last few years have seen many new transport protocols being developed and adopted in the wild. This evolution has been driven by the changing needs of application developers and has been enabled by moving transport services into the application or by tunnelling over an underlying UDP connection.

Application developers are now faced with a genuine choice of different protocols with no clear mechanism for choosing between them. At the same time, the still-limited deployment of some protocols means that the developer must always provide a fall-back to an alternative transport if they want to guarantee the connection will work. This is not a sustainable state of affairs and we believe that in future a new transport API will be needed that provides the mechanisms to facilitate the choice of transport protocol. The first step towards this is to identify the set of Transport Services that a transport protocol is able to expose to the application. We propose doing this in a bottom-up fashion, starting from the list of services available in transport protocols that are specified in RFCs.

7. Contributors and Acknowledgements

Many thanks to the many people that have contributed to this effort so far including Arjuna Sathiaselan, Jon Crowcroft, Marwan Fayed and Bernd Reuther among many others.

D. Ros and M. Welzl were part-funded by the European Community under its Seventh Framework Programme through the Reducing Internet Transport Latency (RITE) project (ICT-317700). T. Moncaster and J. Crowcroft are part-funded by the European Union's Seventh Framework Programme FP7/2007-2013 under the Trilogy 2 project, grant agreement no. 317756.

8. Comments Solicited

To be removed by RFC Editor: This draft is the first step towards an IETF BoF on Transport Services. Comments and questions are encouraged and very welcome. They can be addressed to the current mailing list <transport-services@ifi.uio.no> and/or to the authors. We also have a website at <<https://sites.google.com/site/transportprotocolservices/>>

9. References

9.1. Normative References

- [RFC0768] Postel, J., "User Datagram Protocol", STD 6, RFC 768, August 1980.
- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, September 1981.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3828] Larzon, L-A., Degermark, M., Pink, S., Jonsson, L-E., and G. Fairhurst, "The Lightweight User Datagram Protocol (UDP-Lite)", RFC 3828, July 2004.
- [RFC4340] Kohler, E., Handley, M., and S. Floyd, "Datagram Congestion Control Protocol (DCCP)", RFC 4340, March 2006.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.
- [RFC5246] Dierks, T. and E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.2", RFC 5246, August 2008.

- [RFC6817] Shalunov, S., Hazel, G., Iyengar, J., and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012.
- [RFC6824] Ford, A., Raiciu, C., Handley, M., and O. Bonaventure, "TCP Extensions for Multipath Operation with Multiple Addresses", RFC 6824, January 2013.
- [RFC6951] Tuexen, M. and R. Stewart, "UDP Encapsulation of Stream Control Transmission Protocol (SCTP) Packets for End-Host to End-Host Communication", RFC 6951, May 2013.

9.2. Informative References

- [Bonaventure-Blog]
Bonaventure, O., "Blog Entry: MPTCP used in iOS 7", September 2013.
- [I-D.dreibholz-tsvwg-sctpsocket-multipath]
Dreibholz, T., Becke, M., and H. Adhari, "SCTP Socket API Extensions for Concurrent Multipath Transfer", draft-dreibholz-tsvwg-sctpsocket-multipath-06 (work in progress), July 2013.
- [I-D.ietf-tsvwg-sctp-dtls-encaps]
Tuexen, M., Stewart, R., Jesup, R., and S. Loreto, "DTLS Encapsulation of SCTP Packets", draft-ietf-tsvwg-sctp-dtls-encaps-02 (work in progress), October 2013.
- [I-D.iyengar-minion-protocol]
Jana, J., Cheshire, S., and J. Graessley, "Minion - Wire Protocol", draft-iyengar-minion-protocol-02 (work in progress), October 2013.
- [I-D.thornburgh-adobe-rtmfp]
Thornburgh, M., "Adobe's Secure Real-Time Media Flow Protocol", draft-thornburgh-adobe-rtmfp-10 (work in progress), July 2013.
- [I-D.tuexen-tsvwg-sctp-multipath]
Amer, P., Becke, M., Dreibholz, T., Ekiz, N., Jana, J., Natarajan, P., Stewart, R., and M. Tuexen, "Load Sharing for the Stream Control Transmission Protocol (SCTP)", draft-tuexen-tsvwg-sctp-multipath-07 (work in progress), October 2013.
- [I-D.wing-tsvwg-happy-eyeballs-sctp]

Wing, D. and P. Natarajan, "Happy Eyeballs: Trending Towards Success with SCTP", draft-wing-tsvwg-happy-eyeballs-sctp-02 (work in progress), October 2010.

[PVTCP] Nabi, Z., Moncaster, T., Madhavapeddy, A., Hand, S., and J. Crowcroft, "Evolving TCP: how hard can it be?", Proceedings of ACM CoNEXT 2012, December 2012.

[QUIC] Roskind, J., "Quick UDP Internet Connections", June 2013.

Authors' Addresses

Toby Moncaster (editor)
University of Cambridge
Computer Laboratory
J.J. Thomson Avenue
Cambridge CB3 0FD
UK

Phone: +44 1223 763654
EMail: toby.moncaster@cl.cam.ac.uk

Jon Crowcroft
University of Cambridge
Computer Laboratory
J.J. Thomson Avenue
Cambridge CB3 0FD
UK

Phone: +44 1223 763633
EMail: jon.crowcroft@cl.cam.ac.uk

Michael Welzl
University of Oslo
PO Box 1080 Blindern
Oslo N-0316
Norway

Phone: +47 22 85 24 20
EMail: michawe@ifi.uio.no

David Ros
Telecom Bretagne
Rue de la Chataigneraie, CS 17607
35576 Cesson Sevigne cedex
France

Phone: +33 2 99 12 70 46
EMail: david.ros@telecom-bretagne.eu

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

EMail: tuexen@fh-muenster.de

Network WG
Internet-Draft
Intended status: Standards Track (PS)
Obsoletes: RFC 4594
Updates: RFC 5865
Expires: August 25, 2013

James Polk, ed.
Cisco
Feb, 2013

Standard Configuration of DiffServ Service Classes
draft-polk-tsvwg-rfc4594-update-03.txt

Abstract

This document describes service classes configured with DiffServ and identifies how they are used and how to construct them using Differentiated Services Code Points (DSCPs), traffic conditioners, Per-Hop Behaviors (PHBs), and Active Queue Management (AQM) mechanisms. There is no intrinsic requirement that particular DSCPs, traffic conditioners, PHBs, and AQM be used for a certain service class, but for consistent behavior under the same network conditions, configuring networks as described here is appropriate.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on August 25, 2013.

Copyright Notice

Copyright (c) 2012 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in

Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	3
1.1. Requirements Notation	
1.2. Expected Use in the Network	
1.3. Service Class Definition	
1.4. Key Differentiated Services Concepts	
1.4.1. Queuing	
1.4.1.1. Priority Queuing	
1.4.1.2. Rate Queuing	
1.4.2. Active Queue Management	
1.4.3. Traffic Conditioning	
1.4.4. Differentiated Services Code Point (DSCP)	
1.4.5. Per-Hop Behavior (PHB)	
1.5. Key Service Concepts	
1.5.1. Default Forwarding (DF)	
1.5.2. Assured Forwarding (AF)	
1.5.3. Expedited Forwarding (EF)	1
1.5.4. Class Selector (CS)	1
1.5.5. Admission Control	1
1.6 What Changes are Proposed Here from RFC 4594?.....	1
2. Service Differentiation	1
2.1. Service Classes	1
2.2. Categorization of User Oriented Service Classes	1
2.3. Service Class Characteristics	1
2.4. Service Classes vs. Treatment Aggregate (from RFC 5127)...	2
2.4.1 Examples of Service Classes in Treatment Aggregates...	2
3. Network Control Traffic	2
3.1. Current Practice in the Internet	2
3.2. Network Control Service Class	2
3.3. OAM Service Class	2
4. User Oriented Traffic	3
4.1. Conversational Service Class Group	3
4.1.1 Audio Service Class	3
4.1.2 Video Service Class	3
4.1.3 Hi-Res Service Class	3
4.2. Realtime-Interactive Service Class	3
4.3. Multimedia Conferencing Service Class	3
4.4. Multimedia Streaming Service Class	3
4.5. Broadcast Video Service Class	4
4.6. Low-Latency Data Service Class	4
4.7. Conversational Signaling Service Class	4
4.8. High-Throughput Data Service Class	4
4.9. Standard Service Class	4
4.10. Low-Priority Data	4
5. Additional Information on Service Class Usage	4
5.1. Mapping for NTP	5

5.2. VPN Service Mapping	5
6. Security Considerations	5
7. Contributing Authors	5
8. Acknowledgements	5
9. References	5
9.1. Normative References	5
9.2. Informative References	5
Author's Address	5
Appendix A - Changes	5

1. Introduction

Differentiated Services [RFC2474][RFC2475] provides the ability to mark/label/classify IP packets differently to distinguish how individual packets need to be treated differently through (or throughout) a network on a per hop basis. Local administrators are who configure each router for which Differentiated Services Code Points (DSCP) are to be treated differently, which are to be ignored (i.e., no differentiated treatment), and which DSCPs are to have their packets remarked (to different DSCPs) as they pass through a router. Local administrators are also who assign which applications, or traffic types, should use which DSCPs to receive the treatment the administrators expect within their network.

What most people fail to understand is that DSCPs provide a per hop behavior (PHB) through that router, but not the previous or next router. In this way of understanding PHB markings, one can understand that Differentiated Services (DiffServ) is not a Quality of Service (QoS) mechanism, but rather a Classification of Service (CoS) mechanism.

For instance, there are 64 possible DSCP values, i.e., using 6 bits of the old Type of Service (TOS) byte [RFC0791]. Each can be configured locally to have greater or less treatment relative to any other DSCP with two exceptions*.

- * Expedited Forwarding (EF) [RFC3246] DSCPs have a treatment requirement that any packet marked within an EF class has to be the next packet transmitted out its egress interface. If there are more than one EF marked packet in the queue, obviously the queue sets the order they are transmitted. Further, if there are more than one EF DSCP, local configuration determines if each are treated the same or differently relate to each other EF DSCP. Currently, there are two Expedited Forwarding DSCPs: EF (101110) [RFC3246] and VOICE-ADMIT (101100) [RFC5865].

- * Class Selector 6 (CS6) [RFC2474] is for routing protocol traffic. There are deemed important because if the network does not transmit and receive its routing protocol traffic in a timely manner, the network stops operating properly.

Not all are configured to mean anything other than best effort forwarding by local administrators of a network. Let us say there are 5 DSCPs configured within network A. Network A's administrator chooses and configures which order (obeying the two exceptions noted above) which application packets are treated differently than any other packets within that network (A). The DSCPs are not fixed to a linear order for relative priority on a per hop basis. Further, and this is often the case, there might be packets with the same DSCP arriving at multiple interfaces of a node, each egressing that node out the same interface. At ingress to this node, everything was fine, with no poor behavior or noticeably excessive amount of packets with the same DSCP. However, at the egress interface, there might not be enough capacity to satisfy the load, thus the departing packets transmit at their maximum rate for that DSCP, but have additional latency due to the overload within that one node. This is called fan-in congestion (or problem). By itself, DiffServ will not remedy this problem for the application that is intolerant to added latency because DiffServ only functions within 1 node at a time.

An additional mechanism is needed to ensure each flow or session receives the amount of packets at its destination that the application requires to perform properly; a mechanism such as IntServ, by way of RSVP [RFC2205] or NSIS [RFC4080]. With this added capability to be session aware, something DiffServ is not, the packets transmitted within a single session have a very good probability of arriving in such a way the receiving application can make full use of each. That said, signaling reservations for each session or flow adds complexity, which creates more work for those who maintain and administer such a network. Adding bandwidth and using DiffServ marking is an easier pill to swallow. The deployment of not few, but more and more audio and (particularly bandwidth hogging) video codecs and their respective application rigidity has caused some to conclude that throwing bandwidth at the problem is no longer acceptable.

With this in mind, this document incorporates five of the six new DSCPs from [ID-DSCP] identified as capacity-admitted DSCPs for most of the service classes in this document. As explained in [ID-DSCP], the five new capacity-admitted DSCPs are from Pool 3. [ID-DSCP] goes further to explain that many layer 2 technologies use fewer bits for marking and prioritization. Instead of six bits like DiffServ, they have three bits, which yields a maximum of 8 values, which tend to line up quite well with the TOS field values. Thus, aggregation of DSCPs is typically accomplished by simply ignoring or reducing the number of bits used to the most significant ones available, such as

EF is 101110, at layer 2 this is merely 101;

Broadcast is 011000, at layer 2 this is merely 011.

However, that was not a premise DiffServ was built upon, to merely

reduce the number of bits. In other words, within DiffServ, XXX is not the same as XXX000 (where XXX is the same binary value in both cases).

This document is originally built upon the RFC 4594 effort, while updating some of the usages and expanding the scope for newer applications that are in use today. The idea in RFC 4594 remains true here, to define a set of service classes, each having unique traffic characteristics, and assigning one or more DSCPs to each service class. As much as the focus could be on the DSCP values, it is not. The focus of this document is the unique traffic characteristics of each service class.

There are many services classes defined in this document, not all will be used in each network at any period of time. This consistency packet markings we talk about is for several reasons, including in a network that does not currently implement a certain service class because they do not have that type of traffic in their network, or that the network merely gives that traffic best effort service. Having a solid guideline to know where to progress or reconfigure a network and endpoints to, say from best effort for a particular traffic type, is a very good thing to do more uniformly than not. A fair amount of burden is placed at DS boundaries needing to keep up with which markings turn into which other markings at both ingress and egress to a network. The same holds true for application developers choosing a default DSCP for their application, lacking a guideline means everyone picks for themselves - and usually with a highly inflated sense of self importance for their application or service.

Another point to make is that there are 20+ service classes defined within the IETF, and that is far too many for most service providers to manage effectively. So, they have formed groups around certain aggregation solutions of service classes. One such aggregation group is based on RFC 5127, which defines what it calls a treatment aggregate, which is taking RFC 4594's service classes and placing them each into one of four treatment aggregates for service providers to handle as a group. SG12 within the ITU-T has an alternative that has nine aggregate groups, so there is work to be done to harmonize aggregates of service classes. This discussion is articulated more in section 2.4. At the end of Section 2.4 we have introduced a series of example configurations which provide examples of how only a few service classes - yet still most treatment aggregates - can be configured in example networks.

Does RFC 4594 need updating? That document is an informational guideline on how networks can or should mark certain packet flows with differing traffic characteristics using DiffServ. There are several reasons why this informational RFC lacks the necessary clarity and strength to reach widespread adoption:

- o confusion between RFC 4594 and RFC 5127 [RFC5127], the latter of

which is for aggregating many 6-bit DSCP values into a 3-bit (8 value) field used specifically by service provider (SP) networks.

- o some believe both RFCs are for SPs, while others ignore RFC 5127 and use RFC 4594 as if it were standards track or BCP.
- o some believe RFC 5127 is for SPs only, and want RFC 4594 to reduce the number of DSCPs within its guidelines to recommend using only 3 or 4 DSCPs. This seems to stem from a manageability and operational perspective.
- o some know RFC 4594 is informational and do not follow its guidelines specifically because it is informational.
- o some use DSCP values that are not defined within RFC 4594, making mapping between different networks using similar or identical application flows difficult.
- o some believe enterprise networks should not use either RFC except at the edge of their networks, where they directly connect to SP networks.
- o some argue that the services classes guidance per class is too broad and are therefore not sure in which service class a particular application is to reside.
- o time has shown that video has become a dominant application on the Internet, and many believe it now requires to be treated uniquely in environments that want to. Video also does not always plan nice with audio, so knowing the two use the same transport (RTP) [RFC3550], a means of separation is in order.

Service class definitions are based on the different traffic characteristics and required performance of the applications/services. There are a greater number of service classes in this document than there were when RFC 4594 [RFC4594] was published (the RFC this document intends to obsolete). The required performance of applications/services has also changed since the publication of RFC 4594, specifically in the area of conversational real time communications. As a result, this document has a greater number of real time applications with more granular set of DSCPs due to their different required performances. Like RFC 4594 before, this approach allows those applications with similar traffic characteristics and performance requirements to be placed in the same service class.

The notion of traffic characteristics and required performance is a per application concept, therefore the label name of each service class remains the same on an end-to-end basis, even if we understand that DiffServ is only a PHB and cannot guarantee anything, even packet delivery at the intended destination node. That said, several applications can be configured to have the same DSCP, or

each have different DSCPs that have the same treatment per hop within a network.

Since RFC 4594 was first published, a new concept has been introduced that will appear throughout this document, including DSCP assignments -- the idea of "admitted" traffic, initially introduced into DiffServ within RFC 5865 [RFC5865]. The VOICE-ADMIT Expedited Forwarding class differentiates itself from the EF Expedited Forwarding by having the packets marked be for admitted traffic. This concept of "admitted" traffic is spread throughout the real time traffic classes.

Thus, the document flow is as follows:

- o maintain the general format of RFC 4594;
- o augment the content with the concept of capacity-admission;
- o incorporate more video into this document, as it has become a dominant application in enterprises and other managed networks, as well as on the open public Internet;
- o reduce the discussion on voice and its examples;
- o articulate the subtle differences learned since RFC 4594 was published.

The goal here is to provide a standard configuration for DiffServ DSCP assignments and expected PHBs for enterprises and other managed networks, as well as towards the public Internet with specific traffic characteristics per Service class/DSCP, and example applications shown for each.

This document describes service classes configured with DiffServ and defines how they can be used and how to construct them using Differentiated Services Code Points (DSCPs), and recommends how to construct them using traffic conditioners, Per-Hop Behaviors (PHBs), and Active Queue Management (AQM) mechanisms. There is no intrinsic requirement that particular traffic conditioners, PHBs, and AQM be used for a certain service class, but as a policy and for interoperability it is useful to apply them consistently.

We differentiate services and their characteristics in Section 2. Network control traffic, as well as user oriented traffic are discussed in Sections 3 and 4, respectively. We analyze the security considerations in Section 6. Section 7 offers a tribute to the authors of RFC 4594, from which this document is based. It is in its own section, and not part of the normal acknowledgements portion of each IETF document.

1.1. Requirements Notation

The key words "SHOULD", "SHOULD NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC2119] when they appear in ALL CAPS. These words may also appear in this document in lower case as plain English words, absent their normative meanings.

1.2. Expected Use in the Network

In the Internet today, corporate LANs and ISP WANs are increasingly utilized, to the point in which network congestion is affecting performance of applications. For this reason, congestion, loss, and variation in delay within corporate LANs and ISP backbones is becoming known to the users collectively as "the network is slow for this application" or just "right now" or "for today". Users do not directly detect network congestion. They react to applications that run slow, or to downloads that take too long in their mind(s). The explosion of video traffic on the internet recently has cause much of this, and is often the application the user is using when they have this slowness.

In the past, application slowness occurred for three very good reasons.

- o the networks the user oriented traffic traverses moves through cycles of bandwidth boom and bandwidth bust, the latter of which become apparent with the periodic deployment of new bandwidth-hungry applications.
- o In access networks, the state is often different. This may be because throughput rates are artificially limited or over-subscribed, or because of access network design trade-offs.
- o Other characteristics, such as database design on web servers (that may create contention points, e.g., in filestore) and configuration of firewalls and routers, often look externally like a bandwidth limitation.

The intent of this document is to provide a standardized marking, plus a conditioning and packet treatment strategy so that it can be configured and put into service on any link that is itself congested.

1.3. Service Class Definition

A "service class" represents a similar set of traffic characteristics for delay, loss, and jitter as packets traverse routers in a network. For example, "High-Throughput Data" service class for store-and-forward applications, or a "Broadcast" service

class for minimally time-shifted IPTV or Internet radio broadcasts. Such a service class may be defined locally in a Differentiated Services (DS) domain, or across multiple DS domains, possibly extending end to end. A goal of this document is to have most/all networks assign the same type of traffic the same for consistency.

A service class is a naming convention which is defined as a word, phrase or initialism/acronym representing a set of necessary traffic characteristics of a certain type of data flow. The necessary characteristics of these traffic flows can be realized by the use of defined per-hop behavior that started with [RFC2474]. The actual specification of the expected treatment of a traffic aggregate within a domain may also be defined as a per-domain behavior (PDB) [RFC3086].

Each domain will locally choose to

- o implement one or more service classes with traffic characteristics as defined here, or
- o implement one or more service classes with similar traffic characteristics as defined here, or
- o implement one or more service classes with similar traffic characteristics as defined here and to aggregate one or more service classes to reduce the number of unique DSCPs within their network, or
- o implement one or more non-standard service classes with traffic characteristics not as defined here, or
- o not use DiffServ within their domain.

For example, low delay, low loss, and minimal jitter may be realized using the EF PHB, or with an over-provisioned AF PHB. This must be done with care as it may disrupt the end-to-end performance required by the applications/services. If the packet sizes are similar within an application, but different between two applications, say small voice packets and large video packets, these two applications may not realize optimum results if merged into the same aggregate if there are any bottlenecks in the network. We provide for this flexibility on a per hop or per domain basis within this document.

This document provides standardized markings for traffic with similar characteristics, and usage expectations for PHBs for specific service classes for their consistent implementation.

The Default Forwarding "Standard" service class is REQUIRED; all other service classes are OPTIONAL. That said, each service class lists traffic characteristics that are expected when using that type of traffic. It is RECOMMENDED that applications and protocols that fit a certain traffic characteristic use the appropriate service

class mark, i.e., the DSCP, for consistent behavior. It is expected that network administrators will base their endpoint application and router configuration choices on the level of service differentiation they require to meet the needs of their customers (i.e., their end-users).

1.4. Key Differentiated Services Concepts

In order to fully understand this document, a reader needs to familiarize themselves with the principles of the Differentiated Services Architecture [RFC2474]. We summarize some key concepts here only to provide convenience for the reader, the referenced RFCs providing the authoritative definitions.

1.4.1. Queuing

A queue is a data structure that holds packets that are awaiting transmission. A router interface can only transmit one packet at a time, however fast the interface speed is. If there is only 1 queue at an interface, the packets are transmitted in the order they are received into that queue - called FIFO, or "first in, first out". Sometimes there is a lag in the time between a packets arrives in the queue and when it is transmitted. This delay might be due to lack of bandwidth, or if there are multiple queues on that interface, because a packet is low in priority relative to other packets that are awaiting to transmit. The scheduler is the system entity that chooses which packet is next in line for transmission when more than one packet are awaiting transmission out the same router interface.

1.4.1.1 Priority Queuing

A priority queuing system is a combination of a set of queues and a scheduler that empties the queues (of packets) in priority sequence. When asked for a packet, the scheduler inspects the highest priority queue and, if there is data present, returns a packet from that queue. Failing that, it inspects the next highest priority queue, and so on. A freeway onramp with a stoplight for one lane that allows vehicles in the high-occupancy-vehicle lane to pass is an example of a priority queuing system; the high-occupancy-vehicle lane represents the "queue" having priority.

In a priority queuing system, a packet in the highest priority queue will experience a readily calculated delay. This is proportional to the amount of data remaining to be serialized when the packet arrived plus the volume of the data already queued ahead of it in the same queue. The technical reason for using a priority queue relates exactly to this fact: it limits delay and variations in delay and should be used for traffic that has that requirement.

A priority queue or queuing system needs to avoid starvation of lower-priority queues. This may be achieved through a variety of means, such as admission control, rate control, or network engineering.

1.4.1.2. Rate Queuing

Similarly, a rate-based queuing system is a combination of a set of queues and a scheduler that empties each at a specified rate. An example of a rate-based queuing system is a road intersection with a stoplight. The stoplight acts as a scheduler, giving each lane a certain opportunity to pass traffic through the intersection.

In a rate-based queuing system, such as Weighted Fair Queuing (WFQ) or Weighted Round Robin (WRR), the delay that a packet in any given queue will experience depends on the parameters and occupancy of its queue and the parameters and occupancy of the queues it is competing with. A queue whose traffic arrival rate is much less than the rate at which it lets traffic depart will tend to be empty, and packets in it will experience nominal delays. A queue whose traffic arrival rate approximates or exceeds its departure rate will tend not to be empty, and packets in it will experience greater delay. Such a scheduler can impose a minimum rate, a maximum rate, or both, on any queue it touches.

1.4.2 Active Queue Management

Active Queue Management, or AQM, is a generic name for any of a variety of procedures that use packet dropping or marking to manage the depth of a queue. The canonical example of such a procedure is Random Early Detection (RED), in that a queue is assigned a minimum and maximum threshold, and the queuing algorithm maintains a moving average of the queue depth. While the mean queue depth exceeds the maximum threshold, all arriving traffic is dropped. While the mean queue depth exceeds the minimum threshold but not the maximum threshold, a randomly selected subset of arriving traffic is marked or dropped. This marking or dropping of traffic is intended to communicate with the sending system, causing its congestion avoidance algorithms to kick in. As a result of this behavior, it is reasonable to expect that TCP's cyclic behavior is desynchronized and that the mean queue depth (and therefore delay) should normally approximate the minimum threshold.

A variation of the algorithm is applied in Assured Forwarding PHB [RFC2597], in that the behavior aggregate consists of traffic with multiple DSCP marks, which are intermingled in a common queue. Different minima and maxima are configured for the several DSCPs separately, such that traffic that exceeds a stated rate at ingress is more likely to be dropped or marked than traffic that is within its contracted rate.

1.4.3 Traffic Conditioning

In addition, at the first router in a network that a packet crosses, arriving traffic may be measured and dropped or marked according to a policy, or perhaps shaped on network ingress, as in "A Rate Adaptive Shaper for Differentiated Services" [RFC2963]. This may be used to bias feedback loops, as is done in "Assured Forwarding PHB" [RFC2597], or to limit the amount of traffic in a system, as is done in "Expedited Forwarding PHB" [RFC3246]. Such measurement procedures are collectively referred to as "traffic conditioners". Traffic conditioners are normally built using token bucket meters, for example with a committed rate and burst size, as in Section 1.5.3 of the DiffServ Model [RFC3290]. The Assured Forwarding PHB [RFC2597] uses a variation on a meter with multiple rate and burst size measurements to test and identify multiple levels of conformance.

Multiple rates and burst sizes can be realized using multiple levels of token buckets or more complex token buckets; these are implementation details. The following are some traffic conditioners that may be used in deployment of differentiated services:

- o For Class Selector (CS) PHBs, a single token bucket meter to provide a rate plus burst size control.
- o For Expedited Forwarding (EF) PHB, a single token bucket meter to provide a rate plus burst size control.
- o For Assured Forwarding (AF) PHBs, usually two token bucket meters configured to provide behavior as outlined in "Two Rate Three Color Marker (trTCM)" [RFC2698] or "Single Rate Three Color Marker (srTCM)" [RFC2697]. The two-rate, three-color marker is used to enforce two rates, whereas the single-rate, three-color marker is used to enforce a committed rate with two burst lengths.

1.4.4 Differentiated Services Code Point (DSCP)

The DSCP is a number in the range 0..63 that is placed into an IP packet to mark it according to the class of traffic it belongs in. These are divided into 3 groups, or pools, defined in RFC 2474, arranged as follows:

- o Pool-1 has 32 values designated for standards assignment (of the form 'xxxxx0').
- o Pool-2 has 16 values designated for experimental or local use only (EXP/LU) assignment (of the form 'xxxx11').
- o Pool-3 has 16 values designated for experimental or local use (EXP/LU) assignment (of the form 'xxxx01').

However, pool-3 is allowed to be assigned for one of two reasons,

#1 - if the values in pool-1 are exhausted, or

#2 - if there is a justifiable reason for assigning a pool-3 DSCP prior to pool-1's exhaustion.

1.4.5 Per-Hop Behavior (PHB)

In the end, the mechanisms described above are combined to form a specified set of characteristics for handling different kinds of traffic, depending on the needs of the application. This document seeks to identify useful traffic aggregates and to specify what PHB should be applied to them.

1.5 Key Service Concepts

While Differentiated Services is a general architecture that may be used to implement a variety of services, three fundamental forwarding behaviors have been defined and characterized for general use. These are basic Default Forwarding (DF) behavior for elastic traffic, the Assured Forwarding (AF) behavior, and the Expedited Forwarding (EF) behavior for real-time (inelastic) traffic. The facts that four code points are recommended for AF and that one code point is recommended for EF are arbitrary choices, and the architecture allows any reasonable number of AF and EF classes simultaneously. The choice of four AF classes and one EF class in the current document is also arbitrary, and operators MAY choose to operate more or fewer of either.

The terms "elastic" and "real-time" are defined in [RFC1633], Section 3.1, as a way of understanding broad-brush application requirements. This document should be reviewed to obtain a broad understanding of the issues in quality of service, just as [RFC2475] should be reviewed to understand the data plane architecture used in today's Internet.

1.5.1 Default Forwarding (DF)

The basic forwarding behaviors applied to any class of traffic are those described in [RFC2474] and [RFC2309]. Best-effort service may be summarized as "I will accept your packets" and is typically configured with some bandwidth guarantee. Packets in transit may be lost, reordered, duplicated, or delayed at random. Generally, networks are engineered to limit this behavior, but changing traffic loads can push any network into such a state.

Application traffic in the internet that uses default forwarding is expected to be "elastic" in nature. By this, we mean that the sender of traffic will adjust its transmission rate in response to

changes in available rate, loss, or delay.

For the basic best-effort service, a single DSCP value is provided to identify the traffic, a queue to store it, and active queue management to protect the network from it and to limit delays.

1.5.2 Assured Forwarding (AF)

The Assured Forwarding PHB [RFC2597] behavior is explicitly modeled on Frame Relay's Discard Eligible (DE) flag or ATM's Cell Loss Priority (CLP) capability. It is intended for networks that offer average-rate Service Level Agreements (SLAs) (as FR and ATM networks do). This is an enhanced best-effort service; traffic is expected to be "elastic" in nature. The receiver will detect loss or variation in delay in the network and provide feedback such that the sender adjusts its transmission rate to approximate available capacity.

For such behaviors, multiple DSCP values are provided (two or three, perhaps more using local values) to identify the traffic, a common queue to store the aggregate, and active queue management to protect the network from it and to limit delays. Traffic is metered as it enters the network, and traffic is variously marked depending on the arrival rate of the aggregate. The premise is that it is normal for users occasionally to use more capacity than their contract stipulates, perhaps up to some bound. However, if traffic should be marked or lost to manage the queue, this excess traffic will be marked or lost first.

1.5.3. Expedited Forwarding (EF)

The intent of Expedited Forwarding PHB [RFC3246] is to provide a building block for low-loss, low-delay, and low-jitter services. It can be used to build an enhanced best-effort service: traffic remains subject to loss due to line errors and reordering during routing changes. However, using queuing techniques, the probability of delay or variation in delay is minimized. For this reason, it is generally used to carry voice and for transport of data information that requires "wire like" behavior through the IP network. Voice is an inelastic "real-time" application that sends packets at the rate the codec produces them, regardless of availability of capacity. As such, this service has the potential to disrupt or congest a network if not controlled. It also has the potential for abuse.

To protect the network, at minimum one SHOULD police traffic at various points to ensure that the design of a queue is not overrun, and then the traffic SHOULD be given a low-delay queue (often using priority, although it is asserted that a rate-based queue can do this) to ensure that variation in delay is not an issue, to meet application needs.

1.5.4 Class Selector (CS)

Class Selector, those DSCPs that end in zeros (xxx000), provide support for historical codepoint definitions and PHB requirement. The CS fields provide a limited backward compatibility with legacy practice, as described in [RFC2474], Section 4. Backward compatibility is addressed in two ways,

- First, there are per-hop behaviors that are already in widespread use (e.g., those satisfying the IPv4 Precedence queuing requirements specified in [RFC1812]), and
- this document will continue to permit their use in DS-compliant networks.

In addition, there are some DSCPs that correspond to historical use of the IP Precedence field,

- CS0 (000000) will remain 'Default Forwarding' (also known as 'Best Effort')
- 11xxxx will remain for routing traffic

and will map to PHBs that meet the general requirements specified in [RFC2474], Section 4.2.2.2.

No attempt is made to maintain backward compatibility with the "DTR" or Type of Service (TOS) bits of the IPv4 TOS octet, as defined in [RFC0791] and [RFC1349].

A DS-compliant network can be deployed exclusively by using one or more CS-compliant PHB groups. Thus, for example, codepoint '011000' would map to the same PHB as codepoint '011010'.

1.5.5 Admission Control

Admission control (including refusal when policy thresholds are crossed) can ensure high-quality communication by ensuring the availability of bandwidth to carry a load. Inelastic real-time flows such as Voice over Internet Protocol (VoIP) (audio) or video conferencing services can benefit from use of an admission control mechanism, as generally the audio or video service is configured with over-subscription, meaning that some users may not be able to make a call during peak periods.

For VoIP (audio) service, a common approach is to use signaling protocols such as SIP, H.323, H.248, MEGACO, along with Resource Reservation Protocol (RSVP) to negotiate admittance and use of network transport capabilities. When a user has been authorized to send voice traffic, this admission procedure has verified that data rates will be within the capacity of the network that it will use.

Many RTP voice and video payloads are inelastic and cannot react to loss or delay in any substantive way. For these payload types, the network needs to police at ingress to ensure that the voice traffic stays within its negotiated bounds. Having thus assured a predictable input rate, the network may use a priority queue to ensure nominal delay and variation in delay.

1.5.5.1 Capacity Admitted (*-Admit)

This is a newer group of traffic types that started with RFC 5865 and the Voice-Admit service type. Voice-Admit is an EF class marking but has capacity-admission always applied to it to ensure each of these flows are managed through a network, though not necessarily on an end-to-end basis. This depends on how many networks each flow transits and the load on each transited network. There are a series of new DSCPs proposed in [ID-DSCP], each specifying unique characteristics necessitating a separate marking from what existing before that document.

This document will import in four new '*-Admit' DSCPs from [ID-DSCP], 2 others that are new but not capacity-admitted, one from RFC 5865, and change the existing usage of 2 DSCPs from RFC 4594. This is discussed throughout the rest of this document.

1.6 What Changes are Proposed Here from RFC 4594?

Changing an entire network DiffServ configuration has proven to be a painful experience for both individuals and companies. It is not done very often, and for good reason. This effort is based on experience learned since the publication of RFC 4594 (circa 2006). Audio, once thought to be ok grouped with video, needs to be in separate service classes. Collaboration has taken off, mostly because of mobility, but also because of a worldwide recession that has limited physical travel, and relying on people to do more with their computers. With that in mind, there has been an explosion in application development for the individual (seems everyone has an "app-store"). The following set of bullets has this world - that needs a robust layer 3 - in mind.

- o Scope of document is changed to tighten it up for standards track consideration.
- o This document explicitly states there is a fundamental requirement that a particular DSCP(s) be used for each service class, each with a recommended set of applications to be used by that service class - at least on that individual's externally facing (public) interface.
- o Created the Conversational group of service classes to focus on realtime, mostly bidirectional communications (unless multicast is

used).

- o "Realtime-Interactive"
Moved to (near) realtime TCP-based apps

Why the change? TCP based transports have proven, in certain environments, to be a bidirectional realtime transport, e.g., for multiplayer gaming and virtual desktops applications.

- o "Audio"
Same as Telephony (which is now gone), adds Voice-Admit for capacity-admitted traffic

Why the change? RFC 5865 (Voice-Admit) needed to be added to the Audio service class. Video needed to be separate from audio, hence the name change from Telephony (which includes video) to just audio.

- o "Video"
NEW for video and audio/video conferencing, was in Multimedia-Conferencing service classification

Why the change? Many networks are using the AF4X for video, but others are throwing anything "multimedia" into the same service class (like elastic TCP flows). Video has become so dominant that it should be what mostly goes into one service class.

- o "Hi-Res"
NEW for video and audio/video conferencing

Why the change? This entirely new service class is for local policy based higher end video (think Telepresence). Without congestion, this service class has the same treatment as Video, but if there is any pushback from the network, Hi-Res (note: not married to the name) has a better PHB.

- o "Multimedia-Conferencing"
Now without audio or human video

Why the change? The change is taking bidirectional human audio and video out of this service class. This is all about non-realtime collaboration - even in conjunction with an audio and/or video flow.

- o "Broadcast"
Remains the same, added CS3-Admit for capacity-admitted

Why the change? Removing the "-Video" from the name because there are so many more flows that are Broadcast in realtime than video.

- o "Low-Latency Data"
Remains the same, adds IM & Presence traffic explicitly

Why the change? Merely explicitly stating a place for some

additional traffic types that otherwise could go elsewhere.

- o "Conversational Signaling" (A/V-Sig)
Was 'Signaling'

Why the change? This change is merely a renaming of a service class, and acknowledgement that some of the previous authors inaccurate beliefs that DSCPs were linearly ordered with those values having a higher value definitely getting better treatment than lower values.

2. Service Differentiation

There are practical limits on the level of service differentiation that should be offered in the IP networks. We believe we have defined a practical approach in delivering service differentiation by defining different service classes that networks may choose to support in order to provide the appropriate level of behaviors and performance needed by current and future applications and services. The defined structure for providing services allows several applications having similar traffic characteristics and performance requirements to be grouped into the same service class. This approach provides a lot of flexibility in providing the appropriate level of service differentiation for current and new, yet unknown applications without introducing significant changes to routers or network configurations when a new traffic type is added to the network.

2.1 Service Classes

Traffic flowing in a network can be classified in many different ways. We have chosen to divide it into two groupings, network control and user/subscriber traffic. To provide service differentiation, different service classes are defined in each grouping. The network control traffic group can further be divided into two service classes (see Section 3 for detailed definition of each service class):

- o "Network Control" for routing and network control function.
- o "OAM" (Operations, Administration, and Management) for network configuration and management functions.

The user/subscriber traffic group is broken down into ten service classes to provide service differentiation for all the different types of applications/services (see Section 4 for detailed definition of each service class):

- o Conversational service group consists of three service classes:
 - Audio, which includes both 'admitted' and 'unadmitted' audio

service classes, is for non-one way (i.e., generally bidirectional) audio media packets between human users of smaller size and at a constant delivery rate.

- Hi-Res Video, which includes both 'admitted' and 'unadmitted' Hi-Res Video service classes, is for video traffic from higher end endpoints between human users necessitating different treatment than from desktop or video phone endpoints. This has a clearly business differentiation, and not a technical differentiation - as both Hi-Res-Video and Video will be treated similarly on the wire when no congestion occurs.
- Video, which includes both 'admitted' and 'unadmitted' video service classes, is for video traffic from lower end endpoints between human users necessitating different treatment than from higher end (i.e., Telepresence) endpoints. This has a clearly business differentiation, and not a technical differentiation - as both Hi-Res-Video and Video will be treated similarly on the wire when no congestion occurs.
- o Conversational Signaling service class is for peer-to-peer and client-server signaling and control functions using protocols such as SIP, H.323, H.248, and Media Gateway Control Protocol (MGCP). This traffic needs to not be starved on the network.

Editor's note: RFC 4594 had this DSCP marking as CS5, but with clearly different characteristics (i.e., no sensitivity to jitter or (unreasonable) delay), this DSCP has been moved to a more appropriate (new) value, defined in [ID-DSCP].

- o Real-Time Interactive, which includes both 'admitted' and 'unadmitted' Realtime-Interactive service class, is for bidirectional variable rate inelastic applications that require low jitter and loss and very low delay, such as interactive gaming applications that use RTP/UDP streams for game control commands, and Virtualized Desktop applications between the user and content source, typically in a centralized data center.
- o Multimedia Conferencing, which includes both 'admitted' and 'unadmitted' multimedia conferencing service class, is for applications that require minimal delay, but not like those of realtime application requirements. This service class can be bursty in nature, as well as not transmit packets for some time. Applications such as presentation data or collaborative application sharing will use this service class.
- o Multimedia Streaming, which includes both 'admitted' and 'unadmitted' multimedia streaming service class, is for one-way bufferable streaming media applications such as Video on Demand (VOD) and webcasts.

- o Broadcast, which includes both 'admitted' and 'unadmitted' broadcast service class, is for inelastic streaming media applications that may be of constant or variable rate, requiring low jitter and very low packet loss, such as broadcast TV and live events, video surveillance, and security.
- o Low-Latency Data service class is for data processing applications such as client/server interactions or Instant Messaging (IM) and Presence data.
- o Conversational Signaling (A/V-Sig) service class is for all signaling messages, whether in-band (i.e., along the data path) or out-of-band (separate from the data path), for the purposes of setting up, maintaining, managing and terminating bi- or multi-directional realtime sessions.
- o High-Throughput Data service class is for store and forward applications such as FTP and billing record transfer.
- o Standard service class, commonly called best effort (BE), is for traffic that has not been identified as requiring differentiated treatment.
- o Low-Priority Data service class, which some could call the scavenger class, is for packet flows where bandwidth assurance is not required.

2.2 Categorization of User Oriented Service Classes

The ten defined user/subscriber service classes listed above can be grouped into a small number of application categories. For some application categories, it was felt that more than one service class was needed to provide service differentiation within that category due to the different traffic characteristic of the applications, control function, and the required flow behavior. Figure 1 provides a summary of service class grouping into four application categories.

Application Control Category

- o The Conversational Signaling service class is intended to be used to control applications or user endpoints. Examples of protocols that would use this service class are SIP, XMPP or H.323 for voice and/or video over IP services. User signaling flows have similar performance requirements as Low-Latency Data, they require a separate DSCP to be distinguished other traffic and allow for a treatment that is unique.

Media-Oriented Category

Due to the vast number of new (in process of being deployed) and already-in-use media-oriented services in IP networks, seven service

classes have been defined.

- o Audio service class is intended for Voice-over-IP (VoIP) services. It may also be used for other applications that meet the defined traffic characteristics and performance requirements.
- o Video service class is intended for Video over IP services. It may also be used for other applications that meet the defined traffic characteristics and performance requirements.
- o Hi-Res service class is intended for higher end video services that have the same traffic characteristics as the video service class, but have a business requirement(s) to be treated differently. One example of this is Telepresence video applications.
- o Realtime-Interactive service class is intended for inelastic applications such as desktop virtualization applications and for interactive gaming.
- o Multimedia Conferencing service class is for everything about or within video conferencing solutions that does not include the voice or (human) video components. Several examples are
 - the presentation data part of an IP conference (call).
 - the application sharing part of an IP conference (call).
 - the whiteboarding aspect of an IP conference (call).

Each of the above can be part of a lower end web-conferencing application or part of a higher end Telepresence video conference. Each also has the ability to reduce their transmission rate on detection of congestion. These flows can therefore be classified as rate adaptive and most often more elastic than their voice and video counterparts.

- o Broadcast Video service class is to be used for inelastic traffic flows specifically with minimal buffering expected by the source or destination, which are intended for broadcast HDTV service, as well as for transport of live video (sports or concerts) and audio events.
- o Multimedia Streaming service class is to be used for elastic multimedia traffic flows where buffering is expected. This is the fundamental difference between the Broadcast and multimedia streaming service classes. Multimedia streaming content is typically stored before being transmitted. It is also buffered at the receiving end before being played out. The buffering is sufficiently large to accommodate any variation in transmission rate that is encountered in the network. Multimedia entertainment over IP delivery services that are being developed

can generate both elastic and inelastic traffic flows; therefore, two service classes are defined to address this space, respectively: Multimedia Streaming and Broadcast Video.

Data Category

The data category is divided into three service classes.

- o Low-Latency Data for applications/services that require low delay or latency for bursty but short-lived flows.
- o High-Throughput Data for applications/services that require good throughput for long-lived bursty flows. High Throughput and Multimedia Streaming are close in their traffic flow characteristics with High Throughput being a bit more bursty and not as long-lived as Multimedia Streaming.
- o Low-Priority Data for applications or services that can tolerate short or long interruptions of packet flows. The Low-Priority Data service class can be viewed as "don't care" to some degree.

Best-Effort Category

- o All traffic that is not differentiated in the network falls into this category and is mapped into the Standard service class. If a packet is marked with a DSCP value that is not supported in the network, it SHOULD be forwarded using the Standard service class.

Figure 1, below, provides a grouping of the defined user/subscriber service classes into four categories, with indications of which ones use an independent flow for signaling or control; type of flow behavior (elastic, rate adaptive, or inelastic); and the last column provides end user Class of Service (CoS) rating as defined in ITU-T Recommendation G.1010.

Application Categories	Service Class	Signaled	Flow Behavior	G.1010 Rating
Application Control	A/V Sig	Not applicable	Inelastic	Responsive
Media-	Realtime Interactive	Yes	Inelastic	Interactive
	Audio	Yes	Inelastic	Interactive
	Video	Yes	Inelastic	Interactive
	Hi-Res	Yes	Inelastic	Interactive
	Multimedia	Yes	Rate	Moderately

Oriented	Conferencing		Adaptive	Interactive
	Broadcast	Yes	Inelastic	Responsive
	Multimedia Streaming	Yes	Elastic	Timely
Data	Low-Latency Data	No	Elastic	Responsive
	Conversational Signaling	No	Elastic or Inelastic	Timely
	High-Throughput Data	No	Elastic	Timely
	Low-Priority Data	No	Elastic	Non-critical
Best Effort	Standard	Not Specified		Non-critical

Figure 1. User/Subscriber Service Classes Grouping

Here is a short explanation of the end user CoS category as defined in ITU-T Recommendation G.1010. User oriented traffic is divided into four different categories, namely, interactive, responsive, timely, and non-critical. An example of interactive traffic is between two humans and is most sensitive to delay, loss, and jitter. Another example of interactive traffic is between two servers where very low delay and loss are needed. Responsive traffic is typically between a human and a server but can also be between two servers. Responsive traffic is less affected by jitter and can tolerate longer delays than interactive traffic. Timely traffic is either between servers or servers and humans and the delay tolerance is significantly longer than responsive traffic. Non-critical traffic is normally between servers/machines where delivery may be delay for period of time.

2.3. Service Class Characteristics

This document specifies what network administrators are to expect when configuring service classes identified by their differing characteristics. Figure 2 identifies these service classes along with their characteristics, as well as the tolerance to loss, delay and jitter for each service class. Properly engineered networks to these PHBs will achieve expected results. That said, not all of the identified service classes are expected in each operator's network.

Service Class Name	Traffic Characteristics	Tolerance to		
		Loss	Delay	Jitter
Network Control	Variable size packets, mostly inelastic short messages, but traffic can also burst (BGP)	Low	Low	Yes
Realtime Interactive	Inelastic, mostly variable rate	Low	Very Low	Low
Audio	Fixed-size small packets, inelastic	Very Low	Very Low	Very Low
Video	Fixed-size small-large packets, inelastic	Very Low	Very Low	Very Low
Hi-Res A/V	Fixed-size small-large packets, inelastic	Very Low	Very Low	Very Low
Multimedia Conferencing	Variable size packets, constant transmit interval, rate adaptive, reacts to loss	Low - Medium	Low - Medium	Low - Medium
Multimedia Streaming	Variable size packets, elastic with variable rate	Low - Medium	Medium	High
Broadcast	Constant and variable rate, inelastic, non-bursty flows	Very Low	Medium	Low
Low-Latency Data	Variable rate, bursty short-lived elastic flows	Low	Low - Medium	Yes
Conversational Signaling	Variable size packets, some what bursty short-lived flows	Low	Low	Yes
OAM	Variable size packets, elastic & inelastic flows	Low	Medium	Yes
High-Throughput Data	Variable rate, bursty long-lived elastic flows	Low	Medium - High	Yes
Standard	A bit of everything	Not Specified		
Low-Priority Data	Non-real-time and elastic	High	High	Yes

Figure 2. Service Class Characteristics

Notes for Figure 2: A "Yes" in the jitter-tolerant column implies that received data is buffered at the endpoint and that a moderate level of server or network-induced variation in delay is not expected to affect the application. Applications that use TCP or SCTP as a transport are generally good examples. Routing protocols and peer-to-peer signaling also fall in this class; although loss can create problems in setting up calls, a moderate level of jitter merely makes call placement a little less predictable in duration.

Service classes indicate the required traffic forwarding treatment in order to meet user, application, and/or network expectations. Section 3 defines the service classes that MAY be used for forwarding network control traffic, and Section 4 defines the service classes that MAY be used for forwarding user oriented traffic with examples of intended application types mapped into each service class. Note that the application types are only examples and are not meant to be all-inclusive or prescriptive. Also, note that the service class naming or ordering does not imply any priority ordering. They are simply reference names that are used in this document with associated QoS behaviors that are optimized for the particular application types they support. Network administrators MAY choose to assign different service class names to the service classes that they will support. Figure 3 defines the RECOMMENDED relationship between service classes and DS codepoint assignment with application examples. It is RECOMMENDED that this relationship be preserved end to end.

Service Class Name	DSCP Name	DSCP Value	Application Examples
Network Control	CS6&CS7	11xxxx	Network routing
Realtime Interactive	CS5, CS5-Admit	101000, 101001	Remote/Virtual Desktop and Interactive gaming
Audio	EF Voice-Admit	101110 101100	Voice bearer
Hi-Res A/V	CS4, CS4-Admit	100000, 100001	Conversational Hi-Res Audio/Video bearer
Video	AF41,AF42 AF43	100010,100100 100110	Audio/Video conferencing bearer
Multimedia Conferencing	MC, MC-Admit	011101, 100101	Presentation Data and App Sharing/Whiteboarding
Multimedia Streaming	AF31,AF32 AF33	011010,011100 011110	Streaming video and audio on demand

Broadcast	CS3, CS3-Admit	011000, 011001	Broadcast TV, live events & video surveillance
Low-Latency Data	AF21,AF22 AF23	010010,010100 010110	Client/server trans., Web- based ordering, IM/Pres
Conversational Signaling	A/V-Sig	010001	Conversational signaling
OAM	CS2	010000	OAM&P
High-Throughput Data	AF11,AF12 AF13	001010,001100 001110	Store and forward applications
Low-Priority Data	CS1	001000	Any flow that has no BW assurance
Best Effort	CS0	000000	Undifferentiated applications

Figure 3. DSCP to Service Class Mapping

Notes for Figure 3:

- o Default Forwarding (DF) and Class Selector 0 (CS0) (i.e., Best Effort) provide equivalent behavior and use the same DS codepoint, '000000'.
- o RFC 2474 identifies any DSCP with a value of 11xxxx to be for network control. This remains true, while it removes 12 DSCPs from the overall pool of 64 available DSCP values (the 4 that are x11 from this group are within pool 2 of RFC 2474, and remain as only experimentally assignable/useable).
- o All PHB names that say "-Admit" are to be used only when a capacity-admission protocol is utilized for that or each traffic flow.

Changes from table 3 of RFC 4594 are as follows:

- o The old term "Signaling" was using CS5 (101000), now is exclusively for the "Conversational Signaling" service group using the DSCP name of "A/V-Sig" (010001), which is newly defined in [ID-DSCP]. This is because CS5 aggregates into the 101xxx aggregate when using layer 2 technologies such as 802.3 Ethernet, 802.11 Wireless Ethernet MPLS, etc - each of which only have 3 bits to mark with. A traffic type that can have very large packets and is not delay sensitive (within reason) is not appropriate for have a 101xxx marking. A REQUIRED behavior for this PHB is that it not be starved in any node.

- o "Conversational" is a new term to include all interactive audio and video. The Conversational service group consists of the audio service class, the video service class and the new Hi-Res service class.
- o "Audio" obsoletes the term "Telephony", which has generally not retained the "video" aspect within the IETF, where video is still commonly called out as a separate thing. Audio retains the nonadmitted traffic PHB of EF (101110), while capacity-admitted audio has been added via the RFC 5865 defined PHB Voice-Admit.
- o "Video" now is AF4x, with AF41 specifically for capacity-admitted video traffic, while AF42 and AF43 are nonadmitted video traffic.
- o "Hi-Res A/V", part of the Conversational service group, is created by [ID-DSCP] for an additional business differentiation interactive video marking for higher end traffic. It is within the 100xxx as CS4 (for nonadmitted traffic) and CS4-Admit (100001) (for capacity-admitted traffic).
- o "Realtime Interactive" is now using CS5 (for nonadmitted traffic), but adds a capacity-admitted DSCP CS5-Admit (101001).
- o "Multimedia Conferencing" is no longer using the AF4x DSCPs, rather it will use the new PHB MC (100101) (for capacity-admitted) and MC-Admit (011101) (for nonadmitted traffic).
- o "Multimedia Streaming" retains using AF3x, however, AF31 is now used for capacity-admitted traffic, while AF32/33 are nonadmitted.
- o "Broadcast" replaces "Broadcast Video" using CS3 (for nonadmitted traffic), and adds a capacity-admitted PHB CS3-Admit (011001).

It is expected that network administrators will base their choice of the service classes that they will support on their need.

Figure 4 provides a summary of DiffServ CoS mechanisms that MUST be used for the defined service classes that are further detailed in Sections 3 and 4 of this document. According to what applications/services need to be differentiated, network administrators MAY choose the service class(es) that need to be supported in their network.

Service Class	DSCP	Conditioning at DS Edge	PHB Used	Queuing	AQM
Network Control	CS6/CS7	See Section 3.1	RFC2474	Rate	Yes
Realtime	CS5,	Police using sr+bs	RFC2474	Rate	No

Interactive	CS5- Admit*					
Audio	EF, Voice- Admit*	Police using sr+bs	RFC3246 RFC5865	Priority	No	
Hi-Res A/V	CS4, CS4- Admit*	Police using sr+bs	RFC2474 [ID-DSCP]	Priority	No	
Video	AF41*, AF42 AF43	Using two-rate, three-color marker (such as RFC 2698)	RFC2597	Rate	Yes per DSCP	
Multimedia Conferencing	MC, MC- Admit*	Police using sr+bs	[ID-DSCP] [ID-DSCP]	Rate	No	
Multimedia Streaming	AF31*, AF32 AF33	Using two-rate, three-color marker (such as RFC 2698)	RFC2597	Rate	Yes per DSCP	
Broadcast	CS3, CS3- Admit*	Police using sr+bs	RFC2474 [ID-DSCP]	Rate	No	
Low- Latency Data	AF21 AF22 AF23	Using single-rate, three-color marker (such as RFC 2697)	RFC2597	Rate	Yes per DSCP	
Conversational Signaling	AV-Sig	Police using sr+bs	[ID-DSCP]	Rate	No	
OAM	CS2	Police using sr+bs	RFC2474	Rate	Yes	
High- Throughput Data	AF11 AF12 AF13	Using two-rate, three-color marker (such as RFC 2698)	RFC2597	Rate	Yes per DSCP	
Standard	DF	Not applicable	RFC2474	Rate	Yes	
Low-Priority Data	CS1	Not applicable	RFC3662	Rate	Yes	

Figure 4. Summary of CoS Mechanisms Used for Each Service Class

* denotes each DSCP identified for capacity-admission traffic only.

Notes for Figure 4:

- o Conditioning at DS edge means that traffic conditioning is performed at the edge of the DiffServ network where untrusted user devices are connected to two different administrative DiffServ networks.
- o "sr+bs" represents a policing mechanism that provides single rate with burst size control.
- o The single-rate, three-color marker (srTCM) behavior SHOULD be equivalent to RFC 2697, and the two-rate, three-color marker (trTCM) behavior SHOULD be equivalent to RFC 2698.
- o The PHB for Realtime-Interactive service class SHOULD be configured to provide high bandwidth assurance. It MAY be configured as another EF PHB (one capacity-admitted and one non-capacity-admitted, if both are to be used) that uses relaxed performance parameters and a rate scheduler.
- o The PHB for Multimedia Conferencing service class SHOULD be configured to provide high bandwidth assurance. It MAY be configured as another EF PHB (one capacity-admitted and one non-capacity-admitted, if both are to be used) that uses relaxed performance parameters and a rate scheduler.
- o The PHB for Broadcast service class SHOULD be configured to provide high bandwidth assurance. It MAY be configured as another EF PHB (one capacity-admitted and one non-capacity-admitted, if both are to be used) that uses relaxed performance parameters and a rate scheduler.

2.4. Service Classes vs. Treatment Aggregates (from RFC 5127)

There are misconceptions about the differences between RFC 4594 specified service classes, and RFC 5127 specified treatment aggregates. Often the two are conflated, and more often the phrase service class is used to mean both definitions. Almost all of the text previous to this section is used in defining service classes, and how one service class is different than another service class (based on traffic characteristics of the applications). Treatment aggregates are groupings of service classes with similar, but not identical, traffic characteristics to give similar treatment from a SP's network.

Below is taken from appendix of RFC 5127 as its recommended groupings of service classes into aggregates based in RFC 4594 specified traffic characteristic expectations.

+-----+			
Treatment	Treatment	DSCP	
Aggregate	Aggregate		
	Behavior		

Network Control	CS (RFC 2474)	CS6
Real-Time*	EF (RFC 3246)	EF, CS5, AF41, AF42, AF43, CS4, CS3
Assured Elastic	AF (RFC 2597)	CS2, AF31, AF21, AF11 AF32, AF22, AF12 AF33, AF23, AF13
Elastic	Default (RFC 2474)	Default, (CS0) CS1

Figure 5: RFC 5127 Defined Treatment Aggregate Behavior**

*NOTE: The RFC 5865 created VOICE-ADMIT is absence from the above figure because VOICE-ADMIT was created far later than this recommendation was. VOICE-ADMIT is appropriate for the Realtime Traffic Aggregate.

**NOTE: Figure 5 is directly from the appendix of RFC 5127 as that RFC's recommendation for configuration. This draft does not directly affect RFC 5127. That is left for an update to RFC 5127 itself. Based on the WG's take on this draft, RFC 5127 will necessitate an update to match this document's new service classes and additional DSCPs. The number of treatment aggregates are not expected to change in the RFC 5127 Update draft though, with the possible exception of a new treatment aggregate for capacity admitted flows; meaning there *might* be a 5th treatment aggregate proposed.

Treatment Aggregates are designed to nicely fit into technologies that do not have many different treatment levels to use. Here are 3 examples of technologies limited to an 8-value field,

- MPLS with its 3 Traffic Class (TC) bits [RFC5462].
- IEEE LANs with its 8-value Priority Code Point (PCP) field, as part of the 802.1Q header spec [IEEE1Q].
- IEEE 802.1e, which defines QoS over Wi-Fi, also only defines 8 levels (called User Priority or UP codes) [IEEE1E].

Treatment Aggregates are dependent on service classes to exist. Therefore many service classes can exist without the need to consider the use of treatment aggregates or their 8-value technologies. For example, a Layer 3 VPN can be all that is needed

to transit traffic flows, regardless of desired treatment, between enterprise LAN campuses. From this reality, the number of treatment aggregates has no direct bearing on the number of service classes.

2.4.1 Examples of Service Classes in Treatment Aggregates

It is **not** expected that all traffic characteristics are to be experienced across an SP's network for any given customer. For example, if VOICE-ADMIT is added to the Realtime Treatment Aggregate in Figure 5, there are 8 different service classes within the Realtime Treatment Aggregate. It is not expected that all 8 service classes will be deployed by customer networks traversing SP networks. RFC 5127's Treatment Aggregates are a table to configure which service class goes into which treatment aggregate. If there are 8 services classes in the Realtime treatment aggregate, there is very little difference than if there were one service within that same Realtime treatment aggregate - it would still be necessary to configure that treatment aggregate. Thus, it becomes a question of not

"how many service classes are there that go into treatment aggregates?"

but

"how many treatment aggregates have one or more services classes requiring configuration"?

Of the 4 treatment aggregates shown in Figure 5, if there are existing service classes in only 3 of the aggregates, then only 3 treatment aggregates are necessary. Of the 3 following examples, notice that examples 2 and 3 have the same number of treatment aggregates, but example 3 has more applications in their own service classes.

Examples 2 and 3 are made under the following assumptions:

- this draft's Service Classes and DSCP assignments are utilized.
- the new AF-Sig DSCP in the Assured Elastic treatment aggregate.
- the Audio, Video service classes are in the EF treatment aggregate.
- the VOICE-ADMIT DSCP is in the EF treatment aggregate.

2.4.1.1 Example 1 - Simple Voice Configuration/SLA

For example 1, we have an SP running MPLS and has an SLA to deliver Network Control, Voice and everything else is Best Effort. The

following table would apply to this configuration/SLA:

Applications	Service Class	DSCP(s)	Treatment Aggregate
Network Control	Network Control	CS6	Network Control
Voice	Audio	EF	Realtime
Everything else	DF	Default (CS0)	Elastic

Figure 6. Example 1 Configuration

Insert different treatments for this example
(i.e., AQM, RED, WFQ, colors, etc from above charts)

2.4.1.2 Example 2 - Voice/Video/Surveillance Configuration/SLA

For example 1, we have an SP running MPLS and has an SLA to deliver Control, audio, video, surveillance, audio & video signaling, and everything else is BE

Applications	Service Class	DSCP(s)	Treatment Aggregate
Network Control	Network Control	CS6	Network Control
Voice, video, surveillance	Audio, Video, Broadcast	EF, AF42, CS3	Realtime
audio & video signaling	Conversational Signaling	AV-Sig	Assured Elastic
Everything else	DF	Default (CS0)	Elastic

Figure 7. Example 2 Configuration

Insert different treatments for this example
(i.e., AQM, RED, WFQ, colors, etc from above charts)

2.4.1.2 Example 3 - Complex CAC realtime/Surveillance/+apps Configuration/SLA

For example 1, we have an SP running MPLS and has an SLA to deliver

Control, voice, CAC voice, CAC video, streaming, signaling, LL data, Network Mgmt., and everything else is BE (including non-CAC video because it is not authorized or authenticated on network)

Applications	Service Class	DSCP(s)	Treatment Aggregate
Network Control	Network Control	CS6	Network Control
Voice, CAC-Voice CAC-video, surveillance	Audio, Video, Broadcast	Voice-Admit EF, AF41 CS3	Realtime
audio & video signaling, VOD (streaming), Network Mgmt.	Conversational Signaling, Low- Latency Data, Multimedia Streaming, OAM	AV-Sig AF21 AF31 CS2	Assured Elastic
Everything else	DF	Default (CS0)	Elastic

Figure 8. Example 3 Configuration

Insert different treatments for this example
(i.e., AQM, RED, WFQ, colors, etc from above charts)

3. Network Control Traffic

Network control traffic is defined as packet flows that are essential for stable operation of an administered network, as well as the information exchanged between neighboring networks across a peering point where SLAs are in place. Network control traffic is different from user application control (signaling) that may be generated by some applications or services. Network control traffic is mostly between routers and network nodes (e.g., routing or mgmt protocols) that are used for operating, administering, controlling, or managing whole networks, network parts or just network segments. Network Control Traffic may be split into two service classes, i.e., Network Control and OAM.

3.1. Current Practice in the Internet

Based on today's routing protocols and network control procedures that are used in the Internet, we have determined that CS6 DSCP value SHOULD be used for routing and control and that CS7 DSCP value SHOULD be reserved for future use, specifically if needed for future

routing or control protocols. Network administrators MAY use a Local/Experimental DSCP, any value that contains 11xx11; therefore, they may use a locally defined service class within their network to further differentiate their routing and control traffic.

RECOMMENDED Network Edge Conditioning for CS7 DSCP marked packets:

- o Drop or remark 111xxx packets at ingress to DiffServ network domain.
- o 111xxx marked packets SHOULD NOT be sent across peering points. Exchange of control information across peering points SHOULD be done using CS6 DSCP and the Network Control service class.
- o any internally defined 11xxx1 values, valid within that network domain, be remarked to CS6 upon egress at network peering points.

3.2. Network Control Service Class

The Network Control service class is used for transmitting packets between network devices (routers) that require control (routing) information to be exchanged between similar devices within the administrative domain, as well as across a peering point between adjacent administrative domains. Traffic transmitted in this service class is very important as it keeps the network operational, and it needs to be forwarded in a timely manner.

The Network Control service class SHOULD be configured using the DiffServ CS6 PHB, defined in [RFC2474]. This service class MUST be configured so that the traffic receives a minimum bandwidth guarantee, to ensure that the packets always receive timely service. The configured forwarding resources for Network Control service class MUST be such that the probability of packet drop under peak load is very low. The Network Control service class SHOULD be configured to use a Rate Queuing system such as defined in Section 1.4.1.2 of this document.

The following are examples of protocols and applications that MUST use the Network Control service class if present in a network:

- o Routing packet flows: OSPF, BGP, ISIS, RIP.
- o Control information exchange within and between different administrative domains across a peering point where SLAs are in place.
- o LSP setup using CR-LDP and RSVP-TE.

The following protocols and applications MUST NOT use the Network Control service class:

- o User oriented traffic is not allowed to use this service class.

By user oriented traffic, we mean packet flows that originate from user-controlled end points that are connected to the network.

- o even if originating from a server or a device acting on behalf of a user or endpoint,
- o even if it is application or in-band signaling to establish a connection wholly within a single network or across peering points of/to adjacent networks (e.g., creating a tunnel such as a VPN, or data path control signaling).

The following are traffic characteristics of packet flows in the Network Control service class:

- o Mostly messages sent between routers and network servers.
- o Variable size packets, normally one packet at a time, but traffic can also burst (BGP, OSPF, etc).
- o IGMP, hen is used only for the normal multicast routing purpose.

The REQUIRED DSCP marking is CS6 (Class Selector 6).

RECOMMENDED Network Edge Conditioning:

- o At peering points (between two DiffServ networks) where SLAs are in place, CS6 marked packets MUST be policed, e.g., using a single rate with burst size (sr+bs) token bucket policer to keep the CS6 marked packet flows to within the traffic rate specified in the SLA.
- o CS6 marked packet flows from untrusted sources (for example, end user devices) MUST be dropped or remarked at ingress to the DiffServ network. What a network admin remarks this user oriented traffic to is a matter of local policy, and inspection of the packets can determine which application is used for proper marking to a more appropriate DSCP, such as from table 3. of this document.
- o Packets from users/subscribers are not permitted access to the Network Control service classes.

The fundamental service offered to the Network Control service class is enhanced best-effort service with high bandwidth assurance. Since this service class is used to forward both elastic and inelastic flows, the service SHOULD be engineered so that the Active Queue Management (AQM) [RFC2309] is applied to CS6 marked packets.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth, and the max-threshold specifies the queue depth above which all traffic is dropped or ECN marked. Thus,

in this service class, the following inequality should hold in queue configurations:

- o min-threshold CS6 < max-threshold CS6
- o max-threshold CS6 <= memory assigned to the queue

Note: Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

3.3. OAM Service Class

The OAM (Operations, Administration, and Management) service class is RECOMMENDED for OAM&P (Operations, Administration, and Management and Provisioning) using protocols such as Simple Network Management Protocol (SNMP), Trivial File Transfer Protocol (TFTP), FTP, Telnet, and Common Open Policy Service (COPS). Applications using this service class require a low packet loss but are relatively not sensitive to delay. This service class is configured to provide good packet delivery for intermittent flows.

The OAM service class SHOULD use the Class Selector (CS) PHB defined in [RFC2474]. This service class SHOULD be configured to provide a minimum bandwidth assurance for CS2 marked packets to ensure that they get forwarded. The OAM service class SHOULD be configured to use a Rate Queuing system such as defined in Section 1.4.1.2 of this document.

The following applications SHOULD use the OAM service class:

- o Provisioning and configuration of network elements.
- o Performance monitoring of network elements.
- o Any network operational alarms.

The following are traffic characteristics:

- o Variable size packets.
- o Intermittent traffic flows.
- o Traffic may burst at times.
- o Both elastic and inelastic flows.
- o Traffic not sensitive to delays.

RECOMMENDED DSCP marking:

- o All flows in this service class are marked with CS2 (Class Selector 2).

Applications or IP end points SHOULD pre-mark their packets with CS2 DSCP value. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475].

RECOMMENDED conditioning performed at DiffServ network edge:

- o Packet flow marking (DSCP setting) from untrusted sources (end user devices) SHOULD be verified at ingress to DiffServ network using Multifield (MF) Classification methods, defined in [RFC2475].
- o Packet flows from untrusted sources (end user devices) SHOULD be policed at ingress to DiffServ network, e.g., using single rate with burst size token bucket policer to ensure that the traffic stays within its negotiated or engineered bounds.
- o Packet flows from trusted sources (routers inside administered network) MAY not require policing.
- o Normally OAM&P CS2 marked packet flows are not allowed to flow across peering points. If that is the case, then CS2 marked packets SHOULD be policed (dropped) at both egress and ingress peering interfaces.

The fundamental service offered to "OAM" traffic is enhanced best-effort service with controlled rate. The service SHOULD be engineered so that CS2 marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery. Since this service class is used to forward both elastic and inelastic flows, the service SHOULD be engineered so that Active Queue Management [RFC2309] is applied to CS2 marked packets.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth for each DSCP, and the max-threshold specifies the queue depth above which all traffic with such a DSCP is dropped or ECN marked. Thus, in this service class, the following inequality should hold in queue configurations:

- o min-threshold CS2 < max-threshold CS2
- o max-threshold CS2 <= memory assigned to the queue

Note: Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

4. User Oriented Traffic

User oriented traffic is defined as packet flows between different users or subscribers, or from servers/nodes on behalf of a user. It is the traffic that is sent to or from end-terminals and that

supports a very wide variety of applications and services, to include traffic about a user or application that assists a user communicate. User oriented traffic can be classified in many different ways. What we have articulated throughout this document is a series of non-exhaustive list of categories for classifying user oriented traffic. We differentiated user oriented traffic that is real-time versus non-real-time, elastic or rate-adaptive versus inelastic, sensitive versus insensitive to loss as well as considering whether the traffic is interactive vs. one way communication, its responsiveness, whether it requires timely delivery, and critical versus non-critical. In the final analysis, we used all of the above for service differentiation, mapping application types that seemed to have different sets of performance sensitivities, and requirements to different service classes.

Network administrators can categorize their applications according to the type of behavior that they require and MAY choose to support all or a subset of the defined service classes. At the same time, we include a public facing default DSCP value, with its associated PHB, that is expected for each traffic type to ensure common or pervasive performance. Figure 3 provides some common applications and the forwarding service classes that best support them, based on their performance requirements.

4.1. Conversational Service Class Group

The Conversational Service Class Group consists of 3 different service classes, audio, video, and Hi-Res. We are describing the media sample, or bearer, packets for applications (e.g., RTP from [RFC3550]) that require bi-directional real-time, very low delay, very low jitter, and very low packet loss for relatively constant-rate traffic sources (inelastic traffic sources). It is RECOMMENDED that RTCP feedback use the same service class and be marked with the same DSCP as the bearer traffic for that (audio and/or video) call. This ensures comparable treatment within the network between endpoints.

The signaling to set-up these bearer flows is part of the Conversational Signaling service group that will be discussed later in Section 4. The following 3 subsections will detail what is expected within each bearer service class.

4.1.1 Audio Service Class

This service class MUST be used for IP Audio service.

The fundamental service offered to traffic in the Audio service class is minimum jitter, delay, and packet loss service up to a specified upper bound. There are two PHBs, both EF based, for the Audio service class:

Nonadmitted Audio traffic - MUST use the EF DSCP [RFC3246], and

is for traffic that has not had any capacity admission signaling performed for that flow or session.

Capacity-Admitted Audio traffic - MUST use the Voice-Admit DSCP [RFC5865], and is for traffic that has had any capacity admission signaling performed for that flow or session, e.g., RSVP [RFC2205] or NSIS [RFC4080].

The capacity-admitted Audio traffic operation is similar to an ATM CBR service, which has guaranteed bandwidth and which, if it stays within the negotiated rate, experiences nominal delay and no loss.

The nonadmitted Audio traffic, on the other hand, has had no such explicit guarantee, but has a favorable PHB ensuring high probability of delivery as well as nominal delay and no loss - implicitly assuming there is not too much like marked traffic between users within a flow.

There are two typical scenarios in which audio calls are established, on the public open Internet using protocols such as SIP, XMPP or H.323, or in more managed networks like enterprises or certain service providers which offer a audio service with some feature benefits and take part in the call signaling. These SPs or enterprises also use protocols like SIP, XMPP, H.323, but also use H.248/MEGACO and MGCP.

On the open Internet, typically there is no SP actively involved in the session set-up of calls, and therefore no servers providing assistance or features to help one user contact another user. Often, this traffic is marked or remarked with the DF (i.e., Best Effort) DSCP.

In more managed networks in which one of more operators have active servers aiding the audio call set-up, where DiffServ can be used and preserved to differentiate traffic, networks are offering a service, therefore need to do some, or a lot of engineering to ensure that capacity offered to one or more applications does not exceed the load to the network. Otherwise, the operator will have unhappy users, at least for that application's usage. This is true for any application, but is especially true for inelastic applications in which the application is rigid in its delivery requirements. Audio bearer traffic is typically such an application, video is another such application, but we will get to video in the next subsection.

When a user in a managed network has been authorized to send Audio traffic (i.e., call initiation via the operator's servers was not rejected), the call admission procedure should have verified that the newly admitted flow will be within the capacity of the Audio service class forwarding capability in the network. Capacity verification is a non-trivial thing, and can either be implicitly assumed by the call server(s) based on the operator's network design, or it can be explicitly signaled from an in-data-path

signaling mechanism that verifies the capacity is available now for this call, for each call made within that network. In the latter case, those that do not have verifiable network capacity along the data path are rejected. An in between means method is for call servers to count calls between two or more endpoints. By topologically understanding where the caller and called party is and have configured a known maximum it will allow between the two locations. This is especially true over WAN links that have far less capacity than LAN links or core parts of a network. Network operators will need to understand the topology between any two callers to ensure the appropriate amount of bandwidth is available for an expected number of simultaneous audio calls.

Once more than one bandwidth amount can be used for audio calls, for example - by allowing more than one codec with different bandwidths per codec for such calls, network engineering becomes more difficult. Since the inelastic nature of RTP payloads from this class do not react well to loss or significant delay in any substantive way, the Audio service class MUST forward packets as soon as possible.

The Audio service class that does not have capacity admission performed in the data path MUST use the Expedited Forwarding (EF) PHB, as defined in [RFC3246], so that all packets are forwarded quickly. The Audio service class that does have capacity admission performed in the data path MUST use the Voice-Admit PHB, as defined in [RFC5865], so that all packets are forwarded quickly. The Audio service class SHOULD be configured to use a Priority Queuing system such as that defined in Section 1.4.1.1 of this document.

The following applications SHOULD use the Audio service class:

- o VoIP (G.711, G.729, iLBC and other audio codecs).
- o Voice-band data over IP (modem, fax).
- o T.38 fax over IP.
- o Circuit emulation over IP, virtual wire, etc.
- o IP Virtual Private Network (VPN) service that specifies single-rate, mean network delay that is slightly longer than network propagation delay, very low jitter, and a very low packet loss.

The following are traffic characteristics:

- o Mostly fixed-size packets for VoIP (30, 60, 70, 120 or 200 bytes in size).
- o Packets emitted at constant time intervals.

- o Admission control of new flows is provided by Audio call server, media gateway, gatekeeper, edge router, end terminal, access node or in-data-path signaling that provides flow admission control function.

Applications or IP end points SHOULD pre-mark their packets with EF or Voice-Admit DSCP value, whichever is appropriate. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475].

The RECOMMENDED DSCP marking is EF for nonadmitted audio flows, and Voice-Admit for capacity-admitted flows for the following applications:

- o VoIP (G.711, G.729 and other codecs).
- o Voice-band data over IP (modem and fax).
- o T.38 fax over IP.
- o Circuit emulation over IP, virtual wire, etc.

RECOMMENDED Network Edge Conditioning:

- o Packet flow marking (DSCP setting) from untrusted sources (end user devices) SHOULD be verified at ingress to DiffServ network using Multifield (MF) Classification methods, defined in [RFC2475]. If untrusted, the network edge SHOULD know if capacity-admission has been applied, since the edge router will have taken part in the admission signaling; therefore will know whether EF or Voice-Admit is the proper marking for that flow.
- o Packet flows from untrusted sources (end user devices) SHOULD be policed at ingress to DiffServ network, e.g., using single rate with burst size token bucket policer to ensure that the Audio traffic stays within its negotiated bounds.
- o Policing is OPTIONAL for packet flows from trusted sources whose behavior is ensured via other means (e.g., administrative controls on those systems).
- o Policing of Audio packet flows across peering points where SLA is in place is OPTIONAL as Audio traffic will be controlled by admission control mechanism between peering points.

The fundamental service offered to "Audio" traffic is enhanced best-effort service with controlled rate, very low delay, and very low loss. The service MUST be engineered so that EF marked packet flows have sufficient bandwidth in the network to provide guaranteed delivery. Otherwise, the service will have in place an explicit capacity-admission signaling protocol such as RSVP or NSIS and thus

mark the packets within the flow as Voice-Admit. Normally traffic in this service class does not respond dynamically to packet loss. As such, Active Queue Management [RFC2309] SHOULD NOT be applied to EF marked packet flows.

4.1.2 Video Service Class

The Video service class is for bidirectional applications that require real-time service for both constant and rate-adaptive traffic. SIP and H.323/V2 (and later) versions of video conferencing equipment with constant and dynamic bandwidth adjustment are such applications. The traffic sources in this service class either have a fixed bandwidth requirement (e.g., MPEG2, etc.), or have the ability to dynamically change their transmission rate (e.g., MPEG4/H.264, etc.) based on feedback from the receiver. This feedback SHOULD be accomplished using RTCP [RFC3550]. One approach for this downspeeding has the receiver detect packet loss, thus signaling in an RTCP message to the source the indication of lost (or delayed or out of order) packets in transit. When necessary the source then selects a lower rate encoding codec. When available, the source merely sends less data, resulting in lower resolution of the same visual display.

The Video service class is not for video downloads, webcasts, or single directional video or audio/video traffic of any kind. It is for human-to-human visual interaction between two users, or more if an MTP is used.

Typical video conferencing configurations negotiate the setup of audio/video session using protocols such as SIP and H.323. Just as with networks that have audio traversing them, video typically traverses the same two types of networks: the open big "I" Internet, in which most every type of traffic is best effort (DF), or on a more managed network such as an enterprise or SP's managed network in which servers within either network take part in the call signaling, thereby offering the video service.

When a user in a managed network has been authorized to send video traffic (i.e., call initiation via the operator's servers was not rejected), the call admission procedure should have verified that the newly admitted flow will be within the capacity of the video service class forwarding capability in the network. Capacity verification is a non-trivial thing, and can either be implicitly assumed by the call server(s) based on the operator's network design, or it can be explicitly signaled from an in-data-path signaling mechanism that verifies the capacity is available now for this call, for each call made within that network. In the latter case, those that do not have verifiable network capacity along the data path are rejected. An in between means method is for call servers to count calls between two or more endpoints. By topologically understanding where the caller and called party is and

have configured a known maximum it will allow between the two locations. Video is larger in bandwidth than audio, and the difference can be significant. For example, for a single G.711 audio call that is 80kbps, an associated video bandwidth for the same call can easily be 4Mbps. This is especially true over WAN links that have far less capacity than LAN links or core parts of a network. Network operators will need to understand the topology between any two callers to ensure the appropriate amount of bandwidth is available for an expected number of simultaneous video and/or audio/video calls.

Note that it is OPTIONALLY the case in these networks that the accompanying audio for the video call will be marked as the video is marked (i.e., using the same DSCP), but not always. One reason this has been done is for lip-sync.

The Video service class MUST use the Assured Forwarding (AF) PHB, defined in [RFC2597]. This service class MUST be configured to provide a bandwidth assurance for AF41, AF42, and AF43 marked packets to ensure that they get forwarded. The Video service class SHOULD be configured to use a Rate Queuing system for AF42 and AF43 traffic flows, such as that defined in Section 1.4.1.2 of this document. However, AF41 MUST be designated as the DSCP for use when capacity-admission signaling has been used, such as RSVP or NSIS, to guarantee delivery through the network. AF42 and AF43 will be used for non-admitted video calls, as well as overflows from AF41 sources that send more packets than they have negotiated bandwidth for that call.

The following applications MUST use the Video service class:

- o SIP and H.323/V2 (and later) versions of video conferencing applications (interactive video).
- o Video conferencing applications with rate control or traffic content importance marking.
- o Interactive, time-critical, and mission-critical applications.

NOTE with regards to the above bullet: this usage SHOULD be minimized, else the video traffic will suffer - unless this is engineered into the topology.

The following are traffic characteristics:

- o Variable size packets (i.e., small to large in size).
- o The higher the resolution or change rate between each image, the higher the duration of large packets.
- o Usually constant inter-packet time interval.

- o Can be Variable rate in transmission.
- o Source is capable of reducing its transmission rate based on being told receiver is detecting packet loss (e.g., via RTCP).

Applications or IP end points SHOULD pre-mark their packets with DSCP values as shown below. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475] and mark all packets as AF4x. Note: In this case, the two-rate, three-color marker will be configured to operate in Color-Blind mode.

Mandatory DSCP marking when performed by router closest to source:

- o AF41 = up to specified rate "A", which is dedicated to non-Hi-Res capacity-admitted video traffic.

Note the audio of an A/V call can be marked AF41 as well.

- o AF42 = all non-Hi-Res video traffic marked AF41 in excess of specified rate "A", or new non-admitted video traffic but below specified rate "B".
- o AF43 = in excess of specified rate "B".
- o Where "A" < "B".

Note: One might expect "A" to approximate the peak rates of sum of all admitted video flows, plus the sum of the mean rates and "B" to approximate the sum of the peak rates of those same two flows.

Mandatory DSCP marking when performed by SIP or H.323/V2 videoconferencing equipment:

- o AF41 = SIP or H.323 video conferencing audio stream RTP.
- o AF41 = SIP or H.323 video conferencing video control RTCP.
- o AF41 = SIP or H.323 video conferencing video stream up to specified rate "A".
- o AF42 = SIP or H.323 video conferencing video stream in excess of specified rate "A" but below specified rate "B".
- o AF42 = SIP or H.323 video conferencing video control RTCP, for those video streams that were generated using AF42.
- o AF43 = SIP or H.323 video conferencing video stream in excess of specified rate "B".

- o AF43 = SIP or H.323 video conferencing video control RTCP, for those video streams that were generated using AF43.
- o Where "A" < "B".

Mandatory conditioning performed at DiffServ network edge:

- o The two-rate, three-color marker SHOULD be configured to provide the behavior as defined in trTCM [RFC2698].
- o If packets are marked by trusted sources or a previously trusted DiffServ domain and the color marking is to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Aware mode.
- o If the packet marking is not trusted or the color marking is not to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Blind mode.

The fundamental service offered to nonadmitted "Video" traffic is enhanced best-effort service with controlled rate and delay. The fundamental service offered to capacity-admitted "Video" traffic is a guaranteed service using in-data-path signaling to ensure expected delivery in a timely manner. For a non-admitted video conferencing service, if a 1% packet loss detected at the receiver triggers an encoding rate change, thus dropping to the next lower provisioned video encoding rate then Active Queue Management [RFC2309] SHOULD be used primarily to switch the video encoding rate under congestion, changing from high rate to lower rate, i.e., 1472 kbps to 768 kbps. This rule applies to all AF42 and 43 flows. The probability of loss of AF41 traffic MUST NOT exceed the probability of loss of AF42 traffic, which in turn MUST NOT exceed the probability of loss of AF43 traffic.

Capacity-admitted video service should not result in packet loss. However, administratively this MAY be allowed to cause a purposeful downspeeding event (i.e., a change in resolution or a change in codec) to occur due to congestion.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth for each DSCP, and the max-threshold specifies the queue depth above which all traffic with such a DSCP is dropped or ECN marked. Thus, in this service class, the following inequality should hold in queue configurations:

- o min-threshold AF43 < max-threshold AF43
- o max-threshold AF43 <= min-threshold AF42
- o min-threshold AF42 < max-threshold AF42
- o max-threshold AF42 <= min-threshold AF41

- o min-threshold AF41 < max-threshold AF41
- o max-threshold AF41 <= memory assigned to the queue

Note: This configuration tends to drop AF43 traffic before AF42 and AF42 before AF41. Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

4.1.3 Hi-Res Service Class

The Hi-Res service class is for higher end (i.e., deemed 'more important') bidirectional applications that require real-time service for both constant and rate-adaptive traffic. There are two PHBs, both EF based, for the Hi-Res video conferencing service class:

Nonadmitted Hi-Res traffic - MUST use the CS4 DSCP [RFC2474], and is for traffic that has not had any capacity admission signaling performed for that flow or session.

Capacity-Admitted Hi-Res traffic - MUST use the CS4-Admit DSCP [ID-DSCP], and is for traffic that has had any capacity admission signaling performed for that flow or session, e.g., RSVP [RFC2205] or NSIS [RFC4080].

The capacity-admitted Hi-Res video conferencing traffic operation is similar to an ATM CBR service, which has guaranteed bandwidth and which, if it stays within the negotiated rate, experiences nominal delay and no loss.

SIP and H.323/V2 (and later) versions of video conferencing equipment with constant and dynamic bandwidth adjustment are such applications. The traffic sources in this service class either have a fixed bandwidth requirement (e.g., MPEG2), or have the ability to dynamically change their transmission rate (e.g., MPEG4/H.264) based on feedback from the receiver. This feedback SHOULD be accomplished using RTCP [RFC3550]. One approach for this downspeeding has the receiver detect packet loss, thus signaling in an RTCP message to the source the indication of lost (or delayed or out of order) packets in transit. When necessary the source then selects a lower rate encoding codec. When available, the source merely sends less data, resulting in lower resolution of the same visual display.

The Hi-Res service class, as with the Video service class, is not for video downloads, webcasts, or single directional video or audio/video traffic of any kind. It is for human-to-human visual interaction between two users, or more if a video conference bridge is used.

Typical Hi-Res video conferencing configurations negotiate the setup

of audio/video session using protocols such as SIP and H.323. Hi-Res video conferencing is generally not over the big "I" Internet, rather nearly exclusively over more managed networks such as an enterprise or special purpose SP's managed network in which servers within either network take part in the call signaling, thereby offering the video service. In addition, typically this type of audio/video service has high business expectations for minimized packet loss, pixilation or other issues with the audio/video experience. In the recent past, entire T3s have been dedicated to a signal Hi-Res call; sometimes one T3 per site of a multi-site video conference.

Hi-Res video conferencing often has larger in bandwidth than the typical video call. The audio portion can be increased as well, as stereo capabilities are often necessary to provide an in-room experience from a distance. The difference can be significant (or another step up from just a typical video service). For example, for a single G.711 audio call that is 80kbps, a Hi-Res conference usually runs G.722 wideband audio at 256kbps. Typical video delivery is up to 4Mbps, whereas a Hi-Res conference can have three 1080p/30fps widescreen displays requiring at least 12Mbps, with a burst capability of much more.

If there were no congestion on the wire, the expected treatment between a video service and a Hi-Res conference would be the same. However, it is typically the case that the Hi-Res conferencing flows have more rigid requirements for quality and business-wise, need to be experience far less errors than the regular video service on the same network.

Note that it is likely the case in these networks that the accompanying audio to the Hi-Res video call will be marked as the Hi-Res video is marked (i.e., using the same DSCP).

The Hi-Res service class MUST use the Class Selector 5 (CS4) PHB, defined in [RFC2474], for non-capacity-admitted conferences. While the capacity-admitted Hi-Res conferences MUST use the CS4-Admit PHB, defined in [ID-DSCP]. This service class MUST be configured to provide a bandwidth assurance for CS4 and CS4-Admit marked packets to ensure that they get forwarded. The Hi-Res service class SHOULD be configured to use a Priority Queuing system such as that defined in Section 1.4.1.1 of this document. Further, CS4-Admit will be designated as the DSCP for use when capacity-admission signaling has been used, such as RSVP or NSIS, to guarantee delivery through the network. CS4 will be used for non-admitted Hi-Res conferences, as well as overflows from CS4-Admit sources that send more packets than they have negotiated bandwidth for that call.

The following applications MUST use the Hi-Res service class:

- o SIP and H.323/V2 (and later) versions of Hi-Res video conferencing applications (interactive Hi-Res video).

- o Video conferencing applications with rate control or traffic content importance marking.

The following are traffic characteristics:

- o Variable size packets.
- o The higher the resolution or change rate between each image, the higher the duration of large packets.
- o Usually constant inter-packet time interval.
- o Can be Variable rate in transmission.
- o Source is capable of reducing its transmission rate based on being told receiver is detecting packet loss.

Applications or IP end points SHOULD pre-mark their packets with DSCP values as shown below. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475] and mark all packets as AF4x.

Mandatory DSCP marking when performed by router closest to source:

- o CS4-Admit = up to specified rate "A", which is dedicated to capacity-admitted Hi-Res traffic.

Note the audio of an A/V call can be marked CS4-Admit as well.

- o CS4 = all video traffic marked CS4-Admit in excess of specified rate "A", or new non-admitted video traffic but below specified rate "B".
- o Where "A" < "B".

Note: One might expect "A" to approximate the peak rates of sum of all admitted video flows, plus the sum of the mean rates and "B" to approximate the sum of the peak rates of those same two flows.

Mandatory DSCP marking when performed by SIP or H.323/V2 videoconferencing equipment:

- o CS4-Admit = SIP or H.323 video conferencing audio stream RTP/UDP.
- o CS4-Admit = SIP or H.323 video conferencing video control RTCP/TCP.
- o CS4-Admit = SIP or H.323 video conferencing video stream up to specified rate "A".

- o CS4 = SIP or H.323 video conferencing video stream in excess of specified rate "A" but below specified rate "B".
- o Where "A" < "B".

Mandatory conditioning performed at DiffServ network edge:

- o The two-rate, three-color marker SHOULD be configured to provide the behavior as defined in trTCM [RFC2698].
- o If packets are marked by trusted sources or a previously trusted DiffServ domain and the color marking is to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Aware mode.
- o If the packet marking is not trusted or the color marking is not to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Blind mode.

The fundamental service offered to nonadmitted "Hi-Res" traffic is enhanced best-effort service with controlled rate and delay. The fundamental service offered to capacity-admitted "Hi-Res" traffic is a guaranteed service using in-data-path signaling to ensure expected or timely delivery. Capacity-admitted video service SHOULD NOT result in packet loss. However, administratively this MAY be allowed to cause a purposeful downspeeding event (i.e., a change in resolution or a change in codec) to occur.

4.2. Realtime-Interactive Service Class

The Realtime-Interactive service class is for bidirectional applications that require low loss and jitter and very low delay for constant or variable rate inelastic traffic sources. Interactive gaming applications that do not have the ability to change encoding rates or to mark packets with different importance indications is one good example of such an application. Another set of applications is virtualized desktop applications in which a remote user has a keyboard, mouse and display monitor, but the desktop is virtualized with the memory/processor/applications back in a common data center, requiring near instantaneous feedback on the user's monitor of any changes caused by the application or an action by the user. Rich media protocols for voice and video MUST NOT use the Realtime-Interactive service class, but rather the appropriate service class from the Conversational service group discussed early in Section 4.1.

The Realtime-Interactive service class will use two PHBs:

Nonadmitted Realtime-Interactive traffic - MUST use the CS5 DSCP [RFC2474], and is for traffic that has not had any capacity

admission signaling performed for that flow or session.

Capacity-Admitted Realtime-Interactive traffic - MUST use the CS5-Admit DSCP [ID-DSCP], and is for traffic that has had any capacity admission signaling performed for that flow or session, e.g., RSVP [RFC2205] or NSIS [RFC4080].

The capacity-admitted Realtime-Interactive traffic operation is similar to an ATM CBR service, which has guaranteed bandwidth and which, if it stays within the negotiated rate, experiences nominal delay and no loss.

Either of the above service classes can be configured as EF based by using a relaxed performance parameter and a rate scheduler.

When a user/endpoint has been authorized to start a new session (i.e., joins a networked game or logs onto a virtualized workstation), the admission procedure should have verified that the newly admitted data rates will be within the engineered capacity of the Realtime-Interactive service class. The bandwidth in the core network and the number of simultaneous Realtime-Interactive sessions that can be supported SHOULD be engineered to control traffic load for this service.

This service class SHOULD be configured to provide a high assurance for bandwidth for CS5 PHB, defined in [RFC2474], or CS5-Admit [ID-DSCP] for guaranteed service through a capacity-admission signaling protocol. The Realtime-Interactive service class SHOULD be configured to use a Rate Queuing system such as that defined in Section 1.4.1.2 of this document. Note that either Realtime-Interactive PHB MAY be configured as another EF PHB, specifically CS5-Admit, that uses a relaxed performance parameter and a rate scheduler, in the priority queue as defined in Section 1.4.1.1 of this document.

The following applications MUST use the Realtime-Interactive service class:

- o Interactive gaming and control.
- o Remote Desktop applications
- o Virtualized Desktop applications.
- o Application server-to-application server non-bursty data transfer requiring very low delay.
- o Inelastic, interactive, time-critical, and mission-critical applications requiring very low delay.

The following are traffic characteristics:

- o Variable size packets.
- o Variable rate, though sometimes bursty, which will require engineering of the network to accommodate.
- o Application is sensitive to delay variation between flows and sessions.
- o Lost packets, if any, are usually ignored by application.

RECOMMENDED DSCP marking:

- o All non-admitted flows in this service class are marked with CS5 (Class Selector 5).
- o All capacity-admitted flows in this service class are marked with CS5-Admit.

Applications or IP end points SHOULD pre-mark their packets with CS5 or CS5-Admit DSCP value, depending on whether a capacity-admission signaling protocol is used for a flow. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475].

RECOMMENDED conditioning performed at DiffServ network edge:

- o Packet flow marking (DSCP setting) from untrusted sources (end user devices) SHOULD be verified at ingress to DiffServ network using Multifield (MF) Classification methods defined in [RFC2475].
- o Packet flows from untrusted sources (end user devices) SHOULD be policed at ingress to DiffServ network, e.g., using single rate with burst size token bucket policer to ensure that the traffic stays within its negotiated or engineered bounds.
- o Packet flows from trusted sources (application servers inside administered network) MAY not require policing.
- o Policing of packet flows across peering points MUST adhere to the Service Level Agreement (SLA).

The fundamental service offered to nonadmitted "Realtime-Interactive" traffic is enhanced best-effort service with controlled rate and delay. The fundamental service offered to capacity-admitted "Realtime-Interactive" traffic is a guaranteed service using in-data-path signaling to ensure expected or timely delivery. Capacity-admitted Realtime-Interactive service SHOULD NOT result in packet loss. The service SHOULD be engineered so that CS5 marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery. Normally, traffic in this

service class does not respond dynamically to packet loss. As such, Active Queue Management [RFC2309] SHOULD NOT be applied to CS5 marked packet flows.

4.3. Multimedia Conferencing Service Class

The Multimedia Conferencing service class is for applications that have a low to medium tolerance to delay, and are rate adaptive to lost packets in transit from sources. Presentation Data applications that are operational in conjunction with an audio/video conference is one good example of such an application. Another set of applications is application sharing or whiteboarding applications, also in conjunction to an A/V conference. In either case, the audio & video part of the flow MUST NOT use the Multimedia Conferencing service class, rather the more appropriate service class within the Conversational service group discussed earlier in Section 4.1.

The Multimedia Conferencing service class will use two PHBs:

Nonadmitted Multimedia Conferencing traffic - MUST use the (new) MC DSCP [ID-DSCP], and is for traffic that has not had any capacity admission signaling performed for that flow or session.

Capacity-Admitted Multimedia Conferencing traffic - MUST use the (new) MC-Admit DSCP [ID-DSCP], and is for traffic that has had any capacity admission signaling performed for that flow or session, e.g., RSVP [RFC2205] or NSIS [RFC4080].

The capacity-admitted Multimedia Conferencing traffic operation is similar to an ATM CBR service, which has guaranteed bandwidth and which, if it stays within the negotiated rate, experiences nominal delay and no loss.

When a user/endpoint initiates a presentation data, application sharing or whiteboarding session, it will typically be part of an audio or audio/video conference such as web-conferencing or an existing Telepresence call. The authorization procedure SHOULD be controlled through the coordinated effort to bind the A/V call with the correct Multimedia Conferencing packet flow through some use of identifiers not in scope of this document. The managed network this flow traverse and the number of simultaneous Multimedia Conferencing sessions that can be supported SHOULD be engineered to control traffic load for this service.

The non-capacity admitted Multimedia Conferencing service class SHOULD use the new MC PHB, defined in [ID-DSCP]. This service class SHOULD be configured to provide a high assurance for bandwidth for CS5 marked packets to ensure that they get forwarded. The Multimedia Conferencing service class SHOULD be configured to use a

Rate Queuing system such as that defined in Section 1.4.1.2 of this document. Note that this service class MAY be configured as another EF PHB that uses a relaxed performance parameter, a rate scheduler, and MC-Admit DSCP value, which MUST use the priority queue as defined in Section 1.4.1.1 of this document.

The following applications MUST use the Multimedia Conferencing service class:

- o Presentation Data applications, which can utilize vector graphics, raster graphics or video delivery.
- o Virtualized Desktop applications.
- o Application server-to-application server non-bursty data transfer requiring very low delay.

The following are traffic characteristics:

- o Variable size packets.
- o Variable rate, though sometimes bursty, which will require engineering of the network to accommodate.
- o Application is sensitive to delay variation between flows and sessions.
- o Lost packets, if any, can be ignored by the application.

RECOMMENDED DSCP marking:

- o All non-admitted flows in this service class are marked with the new MC DSCP.
- o All capacity-admitted flows in this service class are marked with MC-Admit.

Applications or IP end points SHOULD pre-mark their packets with the MC DSCP value. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475].

RECOMMENDED conditioning performed at DiffServ network edge:

- o Packet flow marking (DSCP setting) from untrusted sources (end user devices) SHOULD be verified at ingress to DiffServ network using Multifield (MF) Classification methods defined in [RFC2475].
- o Packet flows from untrusted sources (end user devices) SHOULD be policed at ingress to DiffServ network, e.g., using single rate with burst size token bucket policer to ensure that the traffic

stays within its negotiated or engineered bounds.

- o Packet flows from trusted sources (application servers inside administered network) MAY not require policing.
- o Policing of packet flows across peering points MUST adhere to the Service Level Agreement (SLA).

The fundamental service offered to nonadmitted "Multimedia Conferencing" traffic is enhanced best-effort service with controlled rate and delay. The fundamental service offered to capacity-admitted "Multimedia Conferencing" traffic is a guaranteed service using in-data-path signaling to ensure expected or timely delivery. Capacity-admitted Multimedia Conferencing service SHOULD NOT result in packet loss. The service SHOULD be engineered so that Multimedia Conferencing marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery. Normally, traffic in this service class does not respond dynamically to packet loss. As such, Active Queue Management [RFC2309] SHOULD NOT be applied to MC or MC-Admit marked packet flows.

4.4. Multimedia Streaming Service Class

The Multimedia Streaming service class is RECOMMENDED for applications that require near-real-time packet forwarding of variable rate elastic traffic sources that are not as delay sensitive as applications using the Broadcast service class. Such applications include streaming audio and video, some video (movies) on-demand applications, and non-interactive webcasts. In general, the Multimedia Streaming service class assumes that the traffic is buffered at the source/destination; therefore, it is less sensitive to delay and jitter.

The Multimedia Streaming service class MUST use the Assured Forwarding (AF3x) PHB, defined in [RFC2597]. This service class MUST be configured to provide a minimum bandwidth assurance for AF31, AF32, and AF33 marked packets to ensure that they get forwarded. The Multimedia Streaming service class SHOULD be configured to use Rate Queuing system for AF32 and AF33 traffic flows, such as that defined in Section 1.4.1.2 of this document. However, AF31 MUST be designated as the DSCP for use when capacity-admission signaling has been used, such as RSVP or NSIS, to guarantee delivery through the network. AF32 and AF33 will be used for non-admitted streaming flows, as well as overflows from AF31 sources that send more packets than they have negotiated bandwidth for that call.

The following applications SHOULD use the Multimedia Streaming service class:

- o Buffered streaming audio (unicast).

- o Buffered streaming video (unicast).
- o Non-interactive Webcasts.
- o IP VPN service that specifies two rates and is less sensitive to delay and jitter.

The following are traffic characteristics:

- o Variable size packets.
- o The higher the rate, the higher the density of large packets.
- o Variable rate.
- o Elastic flows.
- o Some bursting at start of flow from some applications, as well as an expected stepping up and down on the rate of the flow based on changes in resolution due to network conditions.

Applications or IP end points SHOULD pre-mark their packets with DSCP values as shown below. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475], and mark all packets as AF3x. Note: In this case, the two-rate, three-color marker will be configured to operate in Color-Blind mode.

RECOMMENDED DSCP marking:

- o AF31 = up to specified rate "A".
- o AF32 = all traffic marked AF31 in excess of specified rate "A", or new AF32 traffic but below specified rate "B".
- o AF33 = in excess of specified rate "B".
- o Where "A" < "B".

Note: One might expect "A" to approximate the peak rates of sum of all streaming flows, plus the sum of the mean rates and "B" to approximate the sum of the peak rates of those same two flows.

RECOMMENDED conditioning performed at DiffServ network edge:

- o The two-rate, three-color marker SHOULD be configured to provide the behavior as defined in trTCM [RFC2698].
- o If packets are marked by trusted sources or a previously trusted DiffServ domain and the color marking is to be preserved, then

the two-rate, three-color marker SHOULD be configured to operate in Color-Aware mode.

- o If the packet marking is not trusted or the color marking is not to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Blind mode.

The fundamental service offered to nonadmitted "Multimedia Streaming" traffic is enhanced best-effort service with controlled rate and delay. The fundamental service offered to capacity-admitted "Multimedia Streaming" traffic is a guaranteed service using in-data-path signaling to ensure expected delivery in a reasonable manner. The service SHOULD be engineered so that AF31 marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery. Since the AF3x traffic is elastic and responds dynamically to packet loss, Active Queue Management [RFC2309] SHOULD be used primarily to reduce forwarding rate to the minimum assured rate at congestion points, unless AF31 has had a capacity-admission signaling protocol applied to the flow, such as RSVP or NSIS.

If a capacity-admission signaling protocol applied to the AF31 flow, which SHOULD be the case always, the AF31 PHB MAY be configured as another EF PHB that uses a relaxed performance parameter and a rate scheduler, in the priority queue as defined in Section 1.4.1.1 of this document.

The probability of loss of AF31 traffic MUST NOT exceed the probability of loss of AF32 traffic, which in turn MUST NOT exceed the probability of loss of AF33.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth for each DSCP, and the max-threshold specifies the queue depth above which all traffic with such a DSCP is dropped or ECN marked. Thus, in this service class, the following inequality MUST hold in queue configurations:

- o min-threshold AF33 < max-threshold AF33
- o max-threshold AF33 <= min-threshold AF32
- o min-threshold AF32 < max-threshold AF32
- o max-threshold AF32 <= min-threshold AF31
- o min-threshold AF31 < max-threshold AF31
- o max-threshold AF31 <= memory assigned to the queue

Note#1: this confirmation MUST be modified if AF31 has a capacity-admission signaling protocol applied to those flows, and the above will only apply to AF32 and AF33, while

AF31 (theoretically) has no packet loss.

Note#2: This configuration tends to drop AF33 traffic before AF32 and AF32 before AF31. Note: Many other AQM algorithms exist and are used; they SHOULD be configured to achieve a similar result.

4.5. Broadcast Service Class

The Broadcast service class is RECOMMENDED for applications that require near-real-time packet forwarding with very low packet loss of constant rate and variable rate inelastic traffic sources that are more delay sensitive than applications using the Multimedia Streaming service class. Such applications include broadcast TV, streaming of live audio and video events, some video-on-demand applications, and video surveillance. In general, the Broadcast service class assumes that the destination end point has a dejitter buffer, for video application usually a 2 - 8 video-frame buffer (66 to several hundred of milliseconds), thus expecting far less buffering before play-out than Multimedia Streaming, which can buffer in the seconds to minutes (to hours).

The Broadcast service class will use two PHBs:

Nonadmitted Broadcast traffic - MUST use the CS3 DSCP [RFC2474], and is for traffic that has not had any capacity admission signaling performed for that flow or session.

Capacity-Admitted Broadcast traffic - MUST use the CS3-Admit DSCP [ID-DSCP], and is for traffic that has had any capacity admission signaling performed for that flow or session, e.g., RSVP [RFC2205] or NSIS [RFC4080].

The capacity-admitted Broadcast traffic operation is similar to an ATM CBR service, which has guaranteed bandwidth and which, if it stays within the negotiated rate, experiences nominal delay and no loss.

Either of the above service classes can be configured as EF based by using a relaxed performance parameter and a rate scheduler.

When a user/endpoint initiates a new Broadcast session (i.e., starts an Internet radio application, starts a live Internet A/V event or a camera comes online to do video-surveillance), the admission procedure should be verified within the application that triggers the flow. The newly admitted data rates will SHOULD be within the engineered capacity of the Broadcast service class within that network. The bandwidth in the core network and the number of simultaneous Broadcast sessions that can be supported SHOULD be engineered to control traffic load for this service.

This service class SHOULD be configured to provide high assurance for bandwidth for CS3 marked packets to ensure that they get forwarded. The Broadcast service class SHOULD be configured to use Rate Queuing system such as that defined in Section 1.4.1.2 of this document. Note that either Broadcast PHB MAY be configured as another EF PHB, specifically CS3-Admit, that uses a relaxed performance parameter and a rate scheduler, in the priority queue as defined in Section 1.4.1.1 of this document.

The following applications SHOULD use the Broadcast service class:

- o Video surveillance and security (unicast).
- o TV broadcast including HDTV (likely multicast, but can be unicast).
- o Video on demand (unicast) with control (virtual DVD).
- o Streaming of live audio events (both unicast and multicast).
- o Streaming of live video events (both unicast and multicast).

The following are traffic characteristics:

- o Variable size packets.
- o The higher the rate, the higher the density of large packets.
- o Mixture of variable rate and constant rate flows.
- o Fixed packet emission time intervals.
- o Inelastic flows.

RECOMMENDED DSCP marking:

- o All non-admitted flows in this service class are marked with CS3 (Class Selector 3).
- o All capacity-admitted flows in this service class are marked with CS3-Admit.
- o In some cases, such as those for security and video surveillance applications, it is NOT RECOMMENDED, but allowed to use a different DSCP marking.

If so, then locally user definable (EXP/LU) codepoints in the range '011x11' MAY be used to provide unique traffic identification. The locally administrator definable (EXP/LU, from pool 2 of RFC 2474) codepoint(s) MAY be associated with the PHB that is used for CS3 or CS3-Admit traffic. Furthermore, depending on the network scenario, additional network edge

conditioning policy MAY be needed for the EXP/LU codepoint(s) used.

Applications or IP end points SHOULD pre-mark their packets with CS3 or CS3-Admit DSCP value. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475].

RECOMMENDED conditioning performed at DiffServ network edge:

- o Packet flow marking (DSCP setting) from untrusted sources (end user devices) SHOULD be verified at ingress to DiffServ network using Multifield (MF) Classification methods defined in [RFC2475].
- o Packet flows from untrusted sources (end user devices) SHOULD be policed at ingress to DiffServ network, e.g., using single rate with burst size token bucket policer to ensure that the traffic stays within its negotiated or engineered bounds.
- o Packet flows from trusted sources (application servers inside administered network) MAY not require policing.
- o Policing of packet flows across peering points MUST be performed to the Service Level Agreement (SLA) of those peering entities.

The fundamental service offered to "Broadcast" traffic is enhanced best-effort service with controlled rate and delay. The fundamental service offered to capacity-admitted "Broadcast" traffic is a guaranteed service using in-data-path signaling to ensure expected or timely delivery. Capacity-admitted Broadcast service SHOULD NOT result in packet loss. The service SHOULD be engineered so that CS3 and CS3-Admit marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery. Normally, traffic in this service class does not respond dynamically to packet loss. As such, Active Queue Management [RFC2309] SHOULD NOT be applied to CS3 marked packet flows.

4.6. Low-Latency Data Service Class

The Low-Latency Data service class is RECOMMENDED for elastic and responsive typically client-/server-based applications. Applications forwarded by this service class are those that require a relatively fast response and typically have asymmetrical bandwidth need, i.e., the client typically sends a short message to the server and the server responds with a much larger data flow back to the client. The most common example of this is when a user clicks a hyperlink (~ few dozen bytes) on a web page, resulting in a new web page to be loaded (Kbytes or MBs of data). This service class is configured to provide good response for TCP [RFC1633] short-lived flows that require real-time packet forwarding of variable rate

traffic sources.

The Low-Latency Data service class SHOULD use the Assured Forwarding (AF) PHB, defined in [RFC2597]. This service class SHOULD be configured to provide a minimum bandwidth assurance for AF21, AF22, and AF23 marked packets to ensure that they get forwarded. The Low-Latency Data service class SHOULD be configured to use a Rate Queuing system such as that defined in Section 1.4.1.2 of this document.

The following applications SHOULD use the Low-Latency Data service class:

- o Client/server applications.
- o Systems Network Architecture (SNA) terminal to host transactions (SNA over IP using Data Link Switching (DLSw)).
- o Web-based transactions (E-commerce).
- o Credit card transactions.
- o Financial wire transfers.
- o Enterprise Resource Planning (ERP) applications (e.g., SAP/BaaN).
- o VPN service that supports Committed Information Rate (CIR) with up to two burst sizes.
- o Instant Messaging and Presence protocols (e.g., SIP, XMPP).

The following are traffic characteristics:

- o Variable size packets.
- o Variable packet emission rate.
- o With packet bursts of TCP window size.
- o Short traffic bursts.
- o Source capable of reducing its transmission rate based on detection of packet loss at the receiver or through explicit congestion notification.

Applications or IP end points SHOULD pre-mark their packets with DSCP values as shown below. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475] and mark all packets as AF2x. Note: In this case, the single-rate, three-color marker will be configured to operate in Color-Blind mode.

RECOMMENDED DSCP marking:

- o AF21 = flow stream with packet burst size up to "A" bytes.
- o AF22 = flow stream with packet burst size in excess of "A" but below "B" bytes.
- o AF23 = flow stream with packet burst size in excess of "B" bytes.
- o Where "A" < "B".

RECOMMENDED conditioning performed at DiffServ network edge:

- o The single-rate, three-color marker SHOULD be configured to provide the behavior as defined in srTCM [RFC2697].
- o If packets are marked by trusted sources or a previously trusted DiffServ domain and the color marking is to be preserved, then the single-rate, three-color marker SHOULD be configured to operate in Color-Aware mode.
- o If the packet marking is not trusted or the color marking is not to be preserved, then the single-rate, three-color marker SHOULD be configured to operate in Color-Blind mode.

The fundamental service offered to "Low-Latency Data" traffic is enhanced best-effort service with controlled rate and delay. The service SHOULD be engineered so that AF21 marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery. Since the AF2x traffic is elastic and responds dynamically to packet loss, Active Queue Management [RFC2309] SHOULD be used primarily to control TCP flow rates at congestion points by dropping packets from TCP flows that have large burst size. The probability of loss of AF21 traffic MUST NOT exceed the probability of loss of AF22 traffic, which in turn MUST NOT exceed the probability of loss of AF23. Explicit Congestion Notification (ECN) [RFC3168] MAY also be used with Active Queue Management.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth for each DSCP, and the max-threshold specifies the queue depth above which all traffic with such a DSCP is dropped or ECN marked. Thus, in this service class, the following inequality should hold in queue configurations:

- o min-threshold AF23 < max-threshold AF23
- o max-threshold AF23 <= min-threshold AF22
- o min-threshold AF22 < max-threshold AF22
- o max-threshold AF22 <= min-threshold AF21

- o min-threshold AF21 < max-threshold AF21
- o max-threshold AF21 <= memory assigned to the queue

Note: This configuration tends to drop AF23 traffic before AF22 and AF22 before AF21. Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

4.7. Conversational Signaling Service Class

The Signaling service class is MUST be limited to delay-sensitive signaling traffic only, and then only applying to signaling that involves the Conversational service group. Audio signaling includes signaling between IP phone and soft-switch, soft-client and soft-switch, and media gateway and soft-switch as well as peer-to-peer using various protocols. Video and Hi-Res signaling includes video endpoint to video endpoint, as well as to Media transfer Point (MTP), to call control server(S), etc. This service class is intended to be used for control of voice and video sessions and applications. Protocols using this service class require a relatively fast response, as there are typically several messages of different sizes sent for control of the session. This service class is configured to provide good response for short-lived, intermittent flows that require real-time packet forwarding. This is not the service class for Instant Messaging (IM), that's within the bounds of the Low-Latency Data service class. The Conversational Signaling service class MUST be configured so that the probability of packet drop or significant queuing delay under peak load is very low in IP network segments that provide this interface.

The Conversational Signaling service class MUST use the new A/V-Sig PHB, defined in [ID-DSCP]. This service class MUST be configured to provide a minimum bandwidth assurance for A/V-Sig marked packets to ensure that they get forwarded. In other words, this service class MUST NOT be starved from transmission within a reasonable timeframe, given that the entire Conversational service group depends on these signaling messages successful delivery. Network engineering SHOULD be done to ensure there is roughly 1-4% available per node interface that audio and video traverse. Local conditions MUST be considered when determining exactly how much bandwidth is given to this service class. The Conversational Signaling service class SHOULD be configured to use a Rate Queuing system such as that defined in Section 1.4.1.2 of this document.

The following applications SHOULD use the Conversational Signaling service class:

- o Peer-to-peer IP telephony signaling (e.g., SIP, H.323, XMPP).
- o Peer-to-peer signaling for multimedia applications (e.g., SIP, H.323, XMPP).

- o Peer-to-peer real-time control function.
- o Client-server IP telephony signaling using H.248, MEGACO, MGCP, IP encapsulated ISDN, or other proprietary protocols.
- o Signaling to control IPTV applications using protocols such as IGMP.
- o Signaling flows between high-capacity telephony call servers or soft switches using protocol such as SIP-T. Such high-capacity devices may control thousands of telephony (VoIP) calls.
- o Signaling for one-way video flows, such as RTSP [RFC2326].
- o IGMP, when used for multicast session control such as channel changing in IPTV systems.
- o OPTIONALLY, this service class can be used for on-path reservation signaling for the traffic flows that will use the "admitted" DSCPs. The alternative is to have the on-path signaling (for reservations) use the DSCP within that service class. This provides a similar treatment of the signaling to the data flow, which might be desired.

The following are traffic characteristics:

- o Variable size packets, normally one packet at a time.
- o Intermittent traffic flows.
- o Traffic may burst at times.
- o Delay-sensitive control messages sent between two end points.

RECOMMENDED DSCP marking:

- o All flows in this service class are marked with A/V-Sig.

Applications or IP end points SHOULD pre-mark their packets with A/V-Sig DSCP value. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475].

RECOMMENDED conditioning performed at DiffServ network edge:

- o Packet flow marking (DSCP setting) from untrusted sources (end user devices) SHOULD be verified at ingress to DiffServ network using Multifield (MF) Classification methods defined in [RFC2475].

- o Packet flows from untrusted sources (end user devices) SHOULD be policed at ingress to DiffServ network, e.g., using single rate with burst size token bucket policer to ensure that the traffic stays within its negotiated or engineered bounds.
- o Packet flows from trusted sources (application servers inside administered network) MAY not require policing.
- o Policing of packet flows across peering points in which each peer is participating in the call set-up MUST be performed to the Service Level Agreement (SLA).

The fundamental service offered to "Conversational Signaling" traffic is enhanced best-effort service with controlled rate and delay. The service SHOULD be engineered so that A/V-Sig marked packet flows have sufficient bandwidth in the network to provide high assurance of delivery and low delay. Normally, traffic in this service class does not respond dynamically to packet loss. As such, Active Queue Management [RFC2309] SHOULD NOT be applied to A/V-Sig marked packet flows.

4.8. High-Throughput Data Service Class

The High-Throughput Data service class is RECOMMENDED for elastic applications that require timely packet forwarding of variable rate traffic sources and, more specifically, is configured to provide good throughput for TCP longer-lived flows. TCP [RFC1633] or a transport with a consistent Congestion Avoidance Procedure [RFC2581] [RFC3782] normally will drive as high a data rate as it can obtain over a long period of time. The FTP protocol is a common example, although one cannot definitively say that all FTP transfers are moving data in bulk.

The High-Throughput Data service class SHOULD use the Assured Forwarding (AF) PHB, defined in [RFC2597]. This service class SHOULD be configured to provide a minimum bandwidth assurance for AF11, AF12, and AF13 marked packets to ensure that they are forwarded in a timely manner. The High-Throughput Data service class SHOULD be configured to use a Rate Queuing system such as that defined in Section 1.4.1.2 of this document.

The following applications SHOULD use the High-Throughput Data service class:

- o Store and forward applications.
- o File transfer applications (e.g., FTP, HTTP, etc).
- o Email.
- o VPN service that supports two rates (committed information rate

and excess or peak information rate).

The following are traffic characteristics:

- o Variable size packets.
- o Variable packet emission rate.
- o Variable rate.
- o With packet bursts of TCP window size.
- o Source capable of reducing its transmission rate based on detection of packet loss at the receiver or through explicit congestion notification.

Applications or IP end points SHOULD pre-mark their packets with DSCP values as shown below. If the end point is not capable of setting the DSCP value, then the router topologically closest to the end point SHOULD perform Multifield (MF) Classification, as defined in [RFC2475], and mark all packets as AF1x. Note: In this case, the two-rate, three-color marker will be configured to operate in Color-Blind mode.

RECOMMENDED DSCP marking:

- o AF11 = up to specified rate "A".
- o AF12 = in excess of specified rate "A" but below specified rate "B".
- o AF13 = in excess of specified rate "B".
- o Where "A" < "B".

RECOMMENDED conditioning performed at DiffServ network edge:

- o The two-rate, three-color marker SHOULD be configured to provide the behavior as defined in trTCM [RFC2698].
- o If packets are marked by trusted sources or a previously trusted DiffServ domain and the color marking is to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Aware mode.
- o If the packet marking is not trusted or the color marking is not to be preserved, then the two-rate, three-color marker SHOULD be configured to operate in Color-Blind mode.

The fundamental service offered to "High-Throughput Data" traffic is enhanced best-effort service with a specified minimum rate. The service SHOULD be engineered so that AF11 marked packet flows have

sufficient bandwidth in the network to provide assured delivery. It can be assumed that this class will consume any available bandwidth and that packets traversing congested links may experience higher queuing delays or packet loss. Since the AF_{lx} traffic is elastic and responds dynamically to packet loss, Active Queue Management [RFC2309] SHOULD be used primarily to control TCP flow rates at congestion points by dropping packets from TCP flows that have higher rates first. The probability of loss of AF₁₁ traffic MUST NOT exceed the probability of loss of AF₁₂ traffic, which in turn MUST NOT exceed the probability of loss of AF₁₃. In such a case, if one network customer is driving significant excess and another seeks to use the link, any losses will be experienced by the high-rate user, causing him to reduce his rate. Explicit Congestion Notification (ECN) [RFC3168] MAY also be used with Active Queue Management.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth for each DSCP, and the max-threshold specifies the queue depth above which all traffic with such a DSCP is dropped or ECN marked. Thus, in this service class, the following inequality should hold in queue configurations:

- o min-threshold AF₁₃ < max-threshold AF₁₃
- o max-threshold AF₁₃ <= min-threshold AF₁₂
- o min-threshold AF₁₂ < max-threshold AF₁₂
- o max-threshold AF₁₂ <= min-threshold AF₁₁
- o min-threshold AF₁₁ < max-threshold AF₁₁
- o max-threshold AF₁₁ <= memory assigned to the queue

Note: This configuration tends to drop AF₁₃ traffic before AF₁₂ and AF₁₂ before AF₁₁. Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

4.9. Standard Service Class

The Standard service class is RECOMMENDED for traffic that has not been classified into one of the other supported forwarding service classes in the DiffServ network domain. This service class provides the Internet's "best-effort" forwarding behavior. This service class typically has minimum bandwidth guarantee.

The Standard service class MUST use the Default Forwarding (DF) PHB, defined in [RFC2474], and SHOULD be configured to receive at least a small percentage of forwarding resources as a guaranteed minimum. This service class SHOULD be configured to use a Rate Queuing system such as that defined in Section 1.4.1.2 of this document.

The following applications SHOULD use the Standard service class:

- o Network services, DNS, DHCP, BootP.
- o Any undifferentiated application/packet flow transported through the DiffServ enabled network.

The following is a traffic characteristic:

- o Non-deterministic, mixture of everything.

The RECOMMENDED DSCP marking is DF (Default Forwarding) '000000'.

Network Edge Conditioning:

There is no requirement that conditioning of packet flows be performed for this service class.

The fundamental service offered to the Standard service class is best-effort service with active queue management to limit overall delay. Typical configurations SHOULD use random packet dropping to implement Active Queue Management [RFC2309] or Explicit Congestion Notification [RFC3168], and MAY impose a minimum or maximum rate on the queue.

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth, and the max-threshold specifies the queue depth above which all traffic is dropped or ECN marked. Thus, in this service class, the following inequality should hold in queue configurations:

- o min-threshold DF < max-threshold DF
- o max-threshold DF <= memory assigned to the queue

Note: Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

4.10. Low-Priority Data

The Low-Priority Data service class serves applications that run over TCP [RFC0793] or a transport with consistent congestion avoidance procedures [RFC2581] [RFC3782] and that the user is willing to accept service without guarantees. This service class is specified in [RFC3662] and [QBSS].

The following applications MAY use the Low-Priority Data service class:

- o Any TCP based-application/packet flow transported through the DiffServ enabled network that does not require any bandwidth assurances.

The following is a traffic characteristic:

- o Non-real-time and elastic.

Network Edge Conditioning:

There is no requirement that conditioning of packet flows be performed for this service class.

The RECOMMENDED DSCP marking is CS1 (Class Selector 1).

The fundamental service offered to the Low-Priority Data service class is best-effort service with zero bandwidth assurance. By placing it into a separate queue or class, it may be treated in a manner consistent with a specific Service Level Agreement.

Typical configurations SHOULD use Explicit Congestion Notification [RFC3168] or random loss to implement Active Queue Management [RFC2309].

If RED [RFC2309] is used as an AQM algorithm, the min-threshold specifies a target queue depth, and the max-threshold specifies the queue depth above which all traffic is dropped or ECN marked. Thus, in this service class, the following inequality should hold in queue configurations:

- o min-threshold CS1 < max-threshold CS1
- o max-threshold CS1 <= memory assigned to the queue

Note: Many other AQM algorithms exist and are used; they should be configured to achieve a similar result.

5. Additional Information on Service Class Usage

In this section, we provide additional information on how some specific applications should be configured to use the defined service classes.

5.1. Mapping for NTP

From tests that were performed, indications are that precise time distribution requires a very low packet delay variation (jitter) transport. Therefore, we suggest that the following guidelines for Network Time Protocol (NTP) be used:

- o When NTP is used for providing high-accuracy timing within an administrator's (carrier's) network or to end users/clients, the audio service class SHOULD be used, and NTP packets should be marked with EF DSCP value.

- o For applications that require "wall clock" timing accuracy, the Standard service class should be used, and packets should be marked with DF DSCP.

5.2. VPN Service Mapping

"Differentiated Services and Tunnels" [RFC2983] considers the interaction of DiffServ architecture with IP tunnels of various forms. Further to guidelines provided in RFC 2983, below are additional guidelines for mapping service classes that are supported in one part of the network into a VPN connection. This discussion is limited to VPNs that use DiffServ technology for traffic differentiation.

- o The DSCP value(s) that is/are used to represent a PHB or a PHB group SHOULD be the same for the networks at both ends of the VPN tunnel, unless remarking of DSCP is done as ingress/egress processing function of the tunnel. DSCP marking needs to be preserved along the tunnel, end to end.
- o The VPN MAY be configured to support one or more service classes. It is left up to the administrators of the two networks to agree on the level of traffic differentiation that will be provided in the network that supports VPN service. Service classes are then mapped into the supported VPN traffic forwarding behaviors that meet the traffic characteristics and performance requirements of the encapsulated service classes.
- o The traffic treatment in the network that is providing the VPN service needs to be such that the encapsulated service class or classes receive comparable behavior and performance in terms of delay, jitter, and packet loss and that they are within the limits of the service specified.
- o The DSCP value in the external header of the packet forwarded through the network providing the VPN service can be different from the DSCP value that is used end to end for service differentiation in the end network.
- o The guidelines for aggregation of two or more service classes into a single traffic forwarding treatment in the network that is providing the VPN service is for further study.

6. Security Considerations

This document discusses policy and describes a common policy configuration, for the use of a Differentiated Services Code Point by transports and applications. If implemented as described, it should require that the network do nothing that the network has not already allowed. If that is the case, no new security issues should arise from the use of such a policy.

It is possible for the policy to be applied incorrectly, or for a wrong policy to be applied in the network for the defined service class. In that case, a policy issue exists that the network SHOULD detect, assess, and deal with. This is a known security issue in any network dependent on policy-directed behavior.

A well-known flaw appears when bandwidth is reserved or enabled for a service (for example, voice and/or video transport) and another service or an attacking traffic stream uses it. This possibility is inherent in DiffServ technology, which depends on appropriate packet markings. When bandwidth reservation or a priority queuing system is used in a vulnerable network, the use of authentication and flow admission is recommended. To the author's knowledge, there is no known technical way to respond to an unauthenticated data stream using service that it is not intended to use, and such is the nature of the Internet.

The use of a service class by a user is not an issue when the SLA between the user and the network permits him to use it, or to use it up to a stated rate. In such cases, simple policing is used in the Differentiated Services Architecture. Some service classes, such as Network Control, are not permitted to be used by users at all; such traffic should be dropped or remarked by ingress filters. Where service classes are available under the SLA only to an authenticated user rather than to the entire population of users, authentication and authorization services are required, such as those surveyed in [AUTHMECH].

7. Contributing Authors

This section specifically calls out the authors of RFC 4594, from which this document is based on.

Jozef Babiarez
Nortel Networks

Kwok Ho Chan
Nortel Networks
Email: khchan.work@gmail.com

Fred Baker
Cisco Systems
EMail: fred@cisco.com

Of note, two of the three mentioned authors above worked for Nortel Networks at the time of writing RFC 4594, a company that no longer exists. This author has not seen or heard from those two in many, many years or IETF meetings - as a result of not knowing their new email addresses (or phone numbers).

While much of this document has been rewritten with either edited or

brand new material, there are many short paragraphs that remain as they were from RFC 4594, as well as many sentences that were also left unchanged. Additionally, there were no new graphs, charts, diagrams, or tables introduced, meaning the first 4 tables within this document existed in RFC 4594, created by those authors. Presently, each of those tables contain modified and new information. The last 3 tables, specifically tables 5, 6, & 7 were removed because the examples section was removed.

This author believes there must be proper credit given for all the contributions, including the framework this document retains from that RFC. Periodically, throughout this document, what was written remains the best way of conveying a thought, rule, or otherwise stated behavior or mechanism. Because RFC 4594 was rather large, there is no realistic way of identifying each part that was left untouched. Further, properly quoting that RFC and leaving those sentences embedded in this document would render this document highly unreadable. Another application could be used to show the changes, deletions and additions - but not one that the IETF accepts presently.

This author has created this "Contributing Authors" section as a way of properly identifying those 3 individuals that provided text within this document. We will let the community judge if this is 'good enough' (i.e., rough consensus), or if another way is better.

8. Acknowledgements

The author would like to thank Paul Jones, Glen Lavers, Mo Zanaty, David Benham, Michael Ramalho, Gorrry Fairhurst, David Black, Brian Carpenter, Al Morton, Ruediger Geib and Shitanshu Shah for their comments and questions about this effort that ultimately helped shape this document.

Below are the folks that were acknowledged in RFC 4594, and this author does not want to lose their recognition of contributions to the original effort.

"The authors thank the TSVWG reviewers, David Black, Brian E. Carpenter, and Alan O'Neill for their review and input to this document.

The authors acknowledge a great many inputs, most notably from Bruce Davie, Dave Oran, Ralph Santitoro, Gary Kenward, Francois Audet, Morgan Littlewood, Robert Milne, John Shuler, Nalin Mistry, Al Morton, Mike Pierce, Ed Koehler Jr., Tim Rahrer, Fil Dickinson, Mike Fidler, and Shane Amante. Kimberly King, Joe Zebarth, and Alistair Munroe each did a thorough proofreading,

and the document is better for their contributions."

9. References

9.1. Normative References

- [RFC0791] Postel, J., "Internet Protocol", STD 5, RFC 791, September 1981.
- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, September 1981.
- [RFC1349] Almquist, P., "Type of Service in the Internet Protocol Suite", RFC 1349, July 1992.
- [RFC1812] Baker, F., "Requirements for IP Version 4 Routers", RFC 1812, June 1995.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2309] Braden, B., Clark, D., Crowcroft, J., Davie, B., Deering, S., Estrin, D., Floyd, S., Jacobson, V., Minshall, G., Partridge, C., Peterson, L., Ramakrishnan, K., Shenker, S., Wroclawski, J., and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet", RFC 2309, April 1998.
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, December 1998.
- [RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and W. Weiss, "An Architecture for Differentiated Service", RFC 2475, December 1998.
- [RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", RFC 2597, June 1999.
- [RFC3246] Davie, B., Charny, A., Bennet, J.C., Benson, K., Le Boudec, J., Courtney, W., Davari, S., Firoiu, V., and D. Stiliadis, "An Expedited Forwarding PHB (Per-Hop Behavior)", RFC 3246, March 2002.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3662] Bless, R., Nichols, K., and K. Wehrle, "A Lower Effort Per-Domain Behavior (PDB) for Differentiated Services",

RFC 3662, December 2003.

- [RFC5865] F. Baker, J. Polk, M. Dolly, "A Differentiated Services Code Point (DSCP) for Capacity-Admitted Traffic", RFC 5865, May 2010

9.2. Informative References

- [AUTHMECH] Rescorla, E., "A Survey of Authentication Mechanisms", Work in Progress, September 2005.
- [QBSS] "QBone Scavenger Service (QBSS) Definition", Internet2 Technical Report Proposed Service Definition, March 2001.
- [IEEE1Q] IEEE, 802.1Q Specification
- [IEEE1E] IEEE, 802.1E Wireless LAN User Priority Specification
- [RFC1633] Braden, R., Clark, D., and S. Shenker, "Integrated Services in the Internet Architecture: an Overview", RFC 1633, June 1994.
- [RFC2205] Braden, R., Zhang, L., Berson, S., Herzog, S., and S. Jamin, "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification", RFC 2205, September 1997.
- [RFC2581] Allman, M., Paxson, V., and W. Stevens, "TCP Congestion Control", RFC 2581, April 1999.
- [RFC2697] Heinanen, J. and R. Guerin, "A Single Rate Three Color Marker", RFC 2697, September 1999.
- [RFC2698] Heinanen, J. and R. Guerin, "A Two Rate Three Color Marker", RFC 2698, September 1999.
- [RFC2963] Bonaventure, O. and S. De Cnodder, "A Rate Adaptive Shaper for Differentiated Services", RFC 2963, October 2000.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", RFC 2983, October 2000.
- [RFC2996] Bernet, Y., "Format of the RSVP DCLASS Object", RFC 2996, November 2000.
- [RFC3086] Nichols, K. and B. Carpenter, "Definition of Differentiated Services Per Domain Behaviors and Rules for their Specification", RFC 3086, April 2001.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC

3168, September 2001.

- [RFC3175] Baker, F., Iturralde, C., Le Faucheur, F., and B. Davie, "Aggregation of RSVP for IPv4 and IPv6 Reservations", RFC 3175, September 2001.
- [RFC3290] Bernet, Y., Blake, S., Grossman, D., and A. Smith, "An Informal Management Model for Diffserv Routers", RFC 3290, May 2002.
- [RFC3782] Floyd, S., Henderson, T., and A. Gurtov, "The NewReno Modification to TCP's Fast Recovery Algorithm", RFC 3782, April 2004.
- [RFC5462] L. Andersson, R. Asati, "Multiprotocol Label Switching (MPLS) Label Stack Entry: EXP Field Renamed to Traffic Class Field", RFC 5462, February 2009

Authors' Address

James Polk
3913 Treemont Circle
Colleyville, Texas 76034

Phone: +1.817.271.3552
Email: jmpolk@cisco.com

Appendix A - Changes

Here is a list of all the changes that were captured during the editing process. This will not be a complete list, and others are free to point out what the authors missed, and we'll include that in the next release.

A.1 Since Individual -02 to -03

- o Inserted section 1.6 to explain fundamentally what has changed since RFC 4594, and why changes are necessary.

A.2 Since Individual -01 to -02

- o Added text to the Intro section on the justification from DiffServ Problem Statement draft, as to more of why this update is necessary.
- o Added text to the Intro section expanding on the concept of service classes vs. treatment aggregates (from RFC 5127).

A.3 Since Individual -00 to -01

- o Added Section 2.4 which covers the conflation issues regarding the differences between service classes and treatment aggregates.
- o Added example operational configurations of treatment aggregates applied to this draft's new set of service classes and additional DSCPs.
- o Added references RFC 5865, RFC 5462, IEEE 802.1E and IEEE 802.1Q.

A.4 Since RFC 4594 to Individual Update -00

- o rewrote Intro to emphasize current topics
- o Created a Conversational Service group, comprising the audio, video and Hi-Res service classes - because they have similar characteristics.
- o Incorporated the 6 new DSCPs from [ID-DSCP].
- o moved the example section, en mass, to an appendix that might not be kept for this version. We're not sure it accomplishes what it needs to, and might not provide any real usefulness.
- o Moved 'Realtime-Interactive' service class to CS5, from CS4
- o Changed 'Broadcast Video' service class to 'Broadcast' service class
- o Changed AF4X to 'Video' service class, replacing 'Multimedia Conferencing' service class
- o Moved 'Multimedia Conferencing' service class to different DSCPs
- o Added the 'Hi-Res' service class
- o Removed section 5.1 on signaling choices. It has been included in the main body of the text.
- o Changed document title
- o ...

IETF
Internet-Draft
Intended status: Best Current Practice
Expires: December 14, 2013

G. Shepherd, Ed.
Cisco Systems
June 12, 2013

Multicast UDP Usage Guidelines for Application Designers
draft-shepherd-multicast-udp-guidelines-01

Abstract

The multi-recipient nature of Multicast prevents the use of any point-to-point connection-oriented transport, therefore restricts all Multicast data to be sent over the User Datagram Protocol (UDP). UDP provides a minimal message-passing transport that has no inherent congestion control mechanisms. Because congestion control is critical to the stable operation of the Internet, applications and upper-layer protocols that choose to use Multicast UDP as an Internet service must employ mechanisms to prevent congestion collapse and to establish some degree of fairness with concurrent traffic. This document provides guidelines on the use of UDP for the designers of multicast applications and higher-level protocols.

Status of This Memo

This Internet-Draft is submitted to IETF 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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on December 14, 2013.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of

publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document.

Table of Contents

1. Introduction	2
1.1. Requirements Language	2
2. Multicast UDP Usage Guidelines	3
2.1. Congestion Control Guidelines	3
2.1.1. Bulk Transfer Applications	3
2.1.2. Low Data-Volume Applications	4
2.1.3. UDP Tunnels	4
2.1.4. Message Size Guidelines	4
3. Acknowledgements	5
4. IANA Considerations	5
5. Security Considerations	6
6. References	6
6.1. Normative References	6
6.2. Informative References	6
Appendix A. Additional Stuff	8
Author's Address	8

1. Introduction

The User Datagram Protocol (UDP) [RFC0768] provides a minimal, unreliable, best-effort, message-passing transport to applications and upper-layer protocols (both simply called "applications" in the remainder of this document). [RFC5405] is scoped to provide guidelines for unicast applications only, but all of the general requirements, references, and use cases apply to multicast [RFC1112][RFC4607] UDP application designers as well. This document chooses to only make recommendations in requirements, use cases, and references where they differ from [RFC5405] or are unique for applications sending multicast UDP data (simply called "multicast" in the remainder of this document).

1.1. Requirements Language

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]

2. Multicast UDP Usage Guidelines

2.1. Congestion Control Guidelines

[RFC2309] discusses the dangers of congestion-unresponsive flows and states that "all UDP-based streaming applications should incorporate effective congestion avoidance mechanisms". Many large-scale multicast deployments are within a single administrative domain, and are provisioned over a bandwidth-reserved path or paths where congestion control is less relevant. But there are a growing number of deployment cases where multicast is transiting multiple domains, is tunneled across the unicast Internet, or transits the Internet through a unicast overlay network. This document is only concerned with the latter case of multicast data transiting the larger Internet, either as native IP multicast or encapsulated in a unicast tunnel and does not apply to administratively scoped deployments.

When the multicast traffic exits the administrative domain of a single network or the bi-laterally agreed path between networks, or is tunneled across the unicast Internet either to another multicast network or to an end device, the application SHOULD provide a TCP-compatible aggregate flow over the end-to-end path to each leaf.

There are currently two models of multicast delivery: the Any-Source Multicast (ASM) model as defined in [RFC1112] and the Source-Specific Multicast (SSM) model as defined in [RFC4607]. ASM group members will receive all data sent to the group by any source, while SSM constrains the distribution tree to only one single source. Many congestion-controlled transport protocols are often not applicable to multicast distribution services, or simply won't scale well to very large multicast trees since they require bi-directional communication and adapt the data-rate to accommodate the network conditions to a single receiver. Multicast distribution trees can often fan out to massive numbers of receivers limiting the scalability of an in-band return channel to control the data-rate, and the one-to-many nature of multicast distribution trees prevent adapting the data-rate to individual receiver requirements. For this reason, TCP-compatible aggregate flow for Internet multicast data, either native or tunneled, is the responsibility of the application.

2.1.1. Bulk Transfer Applications

Applications that perform bulk transmission of data over a multicast distribution tree, i.e., applications that exchange more than a small number of UDP datagrams per maximum receiver RTT, SHOULD implement Asynchronous Layered Coding (ALC) [RFC5775], TCP-Friendly Multicast Congestion Control (TFMCC) [RFC4654], Wave and Equation Based Rate Control (WEBRC) [RFC3738], NACK-Oriented Reliable Multicast (NORM)

transport protocol [RFC5740], File Delivery over Unidirectional Transport (FLUTE) [RFC6726], Real Time Protocol/Control Protocol (RTP/RTCP), [RFC3550] or another congestion control scheme following the guidelines of [RFC2887] and utilizing the framework of [RFC3048].

Bulk transfer applications that choose not to implement [RFC4654], [RFC5775], [RFC3738], [RFC5740], [RFC6726], or [RFC3550] SHOULD implement a congestion control scheme that results in bandwidth use that competes fairly with TCP within an order of magnitude. Section 2 of [RFC3551] suggests that applications SHOULD monitor the packet loss rate to ensure that it is within acceptable parameters. Packet loss is considered acceptable if a TCP flow across the same network path under the same network conditions would achieve an average throughput, measured on a reasonable timescale, that is not less than that of the UDP flow. The comparison to TCP cannot be specified exactly, but is intended as an "order-of-magnitude" comparison in timescale and throughput.

Finally, some bulk transfer applications may choose not to implement any congestion control mechanism and instead rely on transmitting across reserved path capacity. This might be an acceptable choice for a subset of restricted networking environments, but is by no means a safe practice for operation in the Internet. When the multicast traffic of such applications leaks out on unprovisioned Internet paths, it can significantly degrade the performance of other traffic sharing the path and even result in congestion collapse. Applications that support an uncontrolled or unadaptive transmission behavior SHOULD NOT do so by default and SHOULD instead require users to explicitly enable this mode of operation.

2.1.2. Low Data-Volume Applications

All of the recommendations in section 3.1.2 of [RFC5405] are applicable to multicast as well.

2.1.3. UDP Tunnels

All of the recommendations in section 3.1.3 of [RFC5405] are applicable to multicast carried inside of unicast UDP tunnels. There are, however deployment cases and solutions where the outer header of a UDP tunnel contains a multicast destination address, such as [RFC6513], but these are primarily deployed in bandwidth reserved environments within a single administrative domain, or between two domains where a bi-laterally agreed upon path and bandwidth is in place and so congestion control is not an issue.

2.1.4. Message Size Guidelines

IP fragmentation lowers the efficiency and reliability of Internet communication. The loss of a single fragment results in the loss of an entire fragmented packet, because even if all other fragments are received correctly, the original packet cannot be reassembled and delivered. This fundamental issue with fragmentation exists for both IPv4 and IPv6, unicast and multicast packets. In addition, some network address translators (NATs) and firewalls drop IP fragments. The network address translation performed by a NAT only operates on complete IP packets, and some firewall policies also require inspection of complete IP packets. Even with these being the case, some NATs and firewalls simply do not implement the necessary reassembly functionality, and instead choose to drop all fragments. Finally, [RFC4963] documents other issues specific to IPv4 fragmentation.

Due to these issues, a multicast application SHOULD NOT send UDP datagrams that result in IP packets that exceed the effective MTU as described in section 3 of [RFC6807]. Consequently, an application SHOULD either use the effective MTU information provided by the Population Count Extensions to Protocol Independent Multicast [RFC6807] or implement path MTU discovery itself [RFC1191][RFC1981][RFC4821] to determine whether the path to each destination will support its desired message size without fragmentation.

If the multicast application is incapable of, or choose not to implement a worst-cast path MTU solution, the application SHOULD assume the maximum MTU of any link will be affected by multiple levels of encapsulation and SHOULD NOT send any packet larger than 1280 bytes.

3. Acknowledgements

This template was derived from an initial version written by Pekka Savola and contributed by him to the xml2rfc project.

This document is part of a plan to make xml2rfc indispensable [DOMINATION].

4. IANA Considerations

This memo includes no request to IANA.

All drafts are required to have an IANA considerations section (see the update of RFC 2434 [I-D.narten-iana-considerations-rfc2434bis] for a guide). If the draft does not require IANA to do anything, the section contains an explicit statement that this is the case (as above). If there are no requirements for IANA, the section will be removed during conversion into an RFC by the RFC Editor.

5. Security Considerations

All drafts are required to have a security considerations section. See RFC 3552 [RFC3552] for a guide.

6. References

6.1. Normative References

- [RFC0768] Postel, J., "User Datagram Protocol", STD 6, RFC 768, August 1980.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [min_ref] authSurName, authInitials., "Minimal Reference", 2006.

6.2. Informative References

- [DOMINATION] Mad Dominators, Inc., "Ultimate Plan for Taking Over the World", 1984, <<http://www.example.com/dominator.html>>.
- [I-D.narten-iana-considerations-rfc2434bis] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", draft-narten-iana-considerations-rfc2434bis-09 (work in progress), March 2008.
- [RFC1112] Deering, S., "Host extensions for IP multicasting", STD 5, RFC 1112, August 1989.
- [RFC1191] Mogul, J. and S. Deering, "Path MTU discovery", RFC 1191, November 1990.
- [RFC1981] McCann, J., Deering, S., and J. Mogul, "Path MTU Discovery for IP version 6", RFC 1981, August 1996.
- [RFC2309] Braden, B., Clark, D., Crowcroft, J., Davie, B., Deering, S., Estrin, D., Floyd, S., Jacobson, V., Minshall, G., Partridge, C., Peterson, L., Ramakrishnan, K., Shenker,

- S., Wroclawski, J., and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet", RFC 2309, April 1998.
- [RFC2629] Rose, M., "Writing I-Ds and RFCs using XML", RFC 2629, June 1999.
- [RFC2887] Handley, M., Floyd, S., Whetten, B., Kermode, R., Vicisano, L., and M. Luby, "The Reliable Multicast Design Space for Bulk Data Transfer", RFC 2887, August 2000.
- [RFC3048] Whetten, B., Vicisano, L., Kermode, R., Handley, M., Floyd, S., and M. Luby, "Reliable Multicast Transport Building Blocks for One-to-Many Bulk-Data Transfer", RFC 3048, January 2001.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3551] Schulzrinne, H. and S. Casner, "RTP Profile for Audio and Video Conferences with Minimal Control", STD 65, RFC 3551, July 2003.
- [RFC3552] Rescorla, E. and B. Korver, "Guidelines for Writing RFC Text on Security Considerations", BCP 72, RFC 3552, July 2003.
- [RFC3738] Luby, M. and V. Goyal, "Wave and Equation Based Rate Control (WEBRC) Building Block", RFC 3738, April 2004.
- [RFC4607] Holbrook, H. and B. Cain, "Source-Specific Multicast for IP", RFC 4607, August 2006.
- [RFC4654] Widmer, J. and M. Handley, "TCP-Friendly Multicast Congestion Control (TFMCC): Protocol Specification", RFC 4654, August 2006.
- [RFC4821] Mathis, M. and J. Heffner, "Packetization Layer Path MTU Discovery", RFC 4821, March 2007.
- [RFC4963] Heffner, J., Mathis, M., and B. Chandler, "IPv4 Reassembly Errors at High Data Rates", RFC 4963, July 2007.
- [RFC5405] Eggert, L. and G. Fairhurst, "Unicast UDP Usage Guidelines for Application Designers", BCP 145, RFC 5405, November 2008.

- [RFC5740] Adamson, B., Bormann, C., Handley, M., and J. Macker, "NACK-Oriented Reliable Multicast (NORM) Transport Protocol", RFC 5740, November 2009.
- [RFC5775] Luby, M., Watson, M., and L. Vicisano, "Asynchronous Layered Coding (ALC) Protocol Instantiation", RFC 5775, April 2010.
- [RFC6513] Rosen, E. and R. Aggarwal, "Multicast in MPLS/BGP IP VPNs", RFC 6513, February 2012.
- [RFC6726] Paila, T., Walsh, R., Luby, M., Roca, V., and R. Lehtonen, "FLUTE - File Delivery over Unidirectional Transport", RFC 6726, November 2012.
- [RFC6807] Farinacci, D., Shepherd, G., Venaas, S., and Y. Cai, "Population Count Extensions to Protocol Independent Multicast (PIM)", RFC 6807, December 2012.

Appendix A. Additional Stuff

This becomes an Appendix.

Author's Address

Greg Shepherd (editor)
Cisco Systems
Tasman Drive
San Jose
USA

Email: gjshep@gmail.com

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: April 23, 2014

R. Stewart
Adara Networks
M. Tuexen
Muenster Univ. of Appl. Sciences
S. Loreto
Ericsson
R. Seggelmann
T-Systems International GmbH
October 20, 2013

A New Data Chunk for Stream Control Transmission Protocol
draft-stewart-tsvwg-sctp-ndata-03.txt

Abstract

The Stream Control Transmission Protocol (SCTP) is a message oriented transport protocol supporting arbitrary large user messages. However, the sender can not interleave different user messages which which causes head of line blocking at the sender side. To overcome this limitation, this document adds a new data chunk to SCTP.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on April 23, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents

carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
2. N-DATA Chunk	3
3. Procedures	4
4. Socket API Considerations	5
5. IANA Considerations	8
6. Security Considerations	9
7. Acknowledgments	9
8. References	9
Authors' Addresses	10

1. Introduction

1.1. Overview

When SCTP [RFC4960] was initially designed it was mainly envisioned for transport of small signaling messages. Late in the design stage it was decided to add support for fragmentation and reassembly of larger messages with the thought that someday Session Initiation Protocol (SIP) [RFC3261] style signaling messages may also need to use SCTP and a single MTU sized message would be too small. Unfortunately this design decision, though valid at the time, did not account for other applications which might send very large messages over SCTP. When such large messages are now sent over SCTP a form of sender side head of line blocking becomes created within the protocol. This head of line blocking is caused by the use of the Transmission Sequence Number (TSN) for two different purposes:

1. As an identifier for DATA chunks to provide a reliable transfer.
2. As an identifier for the sequence of fragments to allow reassembly.

The protocol requires all fragments of a user message to have consecutive TSNs. Therefore the sender can not interleave different messages.

This document describes a new Data chunk called N-DATA. This chunk incorporates all the flags and properties of the current SCTP Data chunk but also adds a new field in its chunk header, the Fragment Sequence Number (FSN). Then the FSN is only used for reassembly and

the TSN only for the reliability. Therefore, the head of line blocking caused by the original design is avoided.

1.2. 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].

2. N-DATA Chunk

The following Figure 1 shows the new data chunk N-DATA.

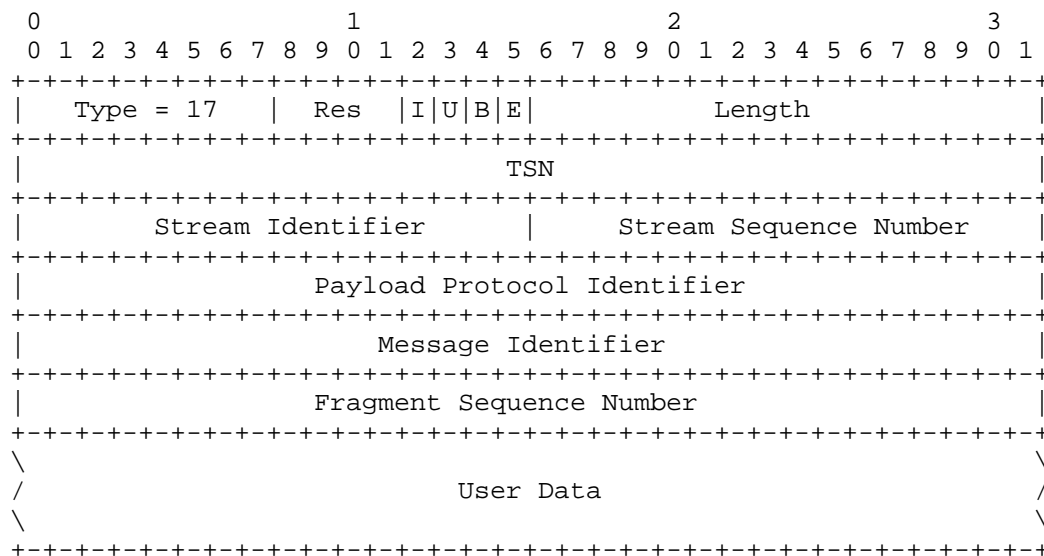


Figure 1: N-DATA chunk format

The only differences between the N-DATA chunk in Figure 1 and the DATA chunk defined in [RFC4960] and [I-D.ietf-tsvwg-sctp-sack-immediately] is the addition of the new Message Identifier (MID) and Fragment Sequence Number (FSN).

Message Identifier (MID): 32 bits (unsigned integer)

The Message Identifier . Please note that the MID is in "network byte order", a.k.a. Big Endian.

Fragment Sequence Number (FSN): 32 bits (unsigned integer)

Identifies the fragment number of this piece of a message. FSN's are unsigned number, the first fragment MUST start at 0 and MUST have the 'B' bit set. The last fragment of a message MUST have

the 'E' bit set. Note that the FSN may wrap completely multiple times allowing arbitrary large messages. Please note that the FSN is in "network byte order", a.k.a. Big Endian.

3. Procedures

3.1. Sender Side Considerations

A sender MUST NOT send a N-DATA chunk unless the peer has indicated its support of the N-DATA chunk type within the Supported Extensions Parameter as defined in [RFC5061].

A sender MUST NOT use the N-DATA chunk unless the user has requested that use via the socket API (see Section 4). This constraint is made since usage of this chunk requires that the application be willing to interleave messages upon reception within an association. This is not the default choice within the socket API (see [RFC6458]) thus the user MUST indicate support to the protocol of the reception of completely interleaved messages. Note that for stacks that do not implement [RFC6458] they may use other methods to indicate interleaved message support and thus enable the usage of the N-DATA chunk, the key is that the the stack MUST know the application has indicated its choice in wanting to use the extension.

Sender side usage of the N-Data chunk is quite simple. Instead of using the TSN for fragmentation purposes, the sender uses the new FSN field to indicate which fragment number is being sent. The first fragment MUST have the 'B' bit set. The last fragment MUST have the 'E' bit set. All other fragments MUST NOT have the 'B' or 'E' bit set. If the 'I' bit is set the 'E' bit MUST also be set, i.e. the 'I' bit may only be set on the last fragment of a message. All other properties of the existing SCTP DATA chunk also apply to the N-DATA chunk, i.e. congestion control as well as receiver window conditions MUST be observed as defined in [RFC4960].

Note that the usage of this chunk should also imply late binding of the actual TSN to any chunk being sent. This way other messages from other streams may be interleaved with the fragmented message.

The sender MUST NOT have more than one ordered fragmented message being produced in any one stream. The sender MUST NOT have more than one un-ordered fragmented message being produced in any one stream. The sender MAY have one ordered and one unordered fragmented message being produced within a single stream. At any time multiple streams MAY be producing an ordered or unordered fragmented message.

3.2. Receiver Side Considerations

Upon reception of an SCTP packet containing a N-DATA chunk if the message needs to be reassembled, then the receiver MUST use the FSN for reassembly of the message and not the TSN. Note that a non-fragmented messages is indicated by the fact that both the 'E' and 'B' bits are set. An ordered or unordered fragmented message is thus identified with any message not having both bits set.

4. Socket API Considerations

This section describes how the socket API defined in [RFC6458] is extended to allow applications to use the extension described in this document.

Please note that this section is informational only.

4.1. Socket Options

option name	data type	get	set
SCTP_NDATA_ENABLE	int	X	X
SCTP_PLUGGABLE_SS	struct sctp_assoc_value	X	X
SCTP_SS_VALUE	struct sctp_stream_value	X	X

4.1.1. Enable or Disable the Interleaving Capability (SCTP_NDATA_ENABLE)

A new socket option to turn on/off the usage of the N-DATA chunk. Turning this option on only effect future associations, and MUST be turned on for the protocol stack to indicate support of the N-DATA chunk to the peer during association setup. Turning this option off, will prevent the N-DATA chunk from being indicated supported in future associations, and will also prevent current associations from producing N-DATA chunks for future large fragmented messages. Note that this does not stop the peer from sending N-DATA chunks.

An N-DATA chunk aware application should also set the fragment interleave level to 2. This allows the reception from multiple streams simultaneously. Failure to set this option can possibly lead to application deadlock.

4.1.2. Get or Set the Stream Scheduler (SCTP_PLUGGABLE_SS)

A stream scheduler can be selected with the SCTP_PLUGGABLE_SS option for `setsockopt()`. The struct `sctp_assoc_value` is used to specify the association for which the scheduler should be changed and the value of the desired algorithm.

The definition of struct `sctp_assoc_value` is the same as in [RFC6458]:

```
struct sctp_assoc_value {
    sctp_assoc_t assoc_id;
    uint32_t assoc_value;
};
```

`assoc_id`: Holds the identifier for the association of which the scheduler should be changed. The special `SCTP_{FUTURE|CURRENT|ALL}_ASSOC` can also be used. This parameter is ignored for one-to-one style sockets.

`assoc_value`: This specifies which scheduler is used. The following constants can be used:

`SCTP_SS_DEFAULT`: The default scheduler used by the SCTP implementation. Typical values are `SCTP_SS_ROUND_ROBIN` or `SCTP_SS_FIRST_COME`.

`SCTP_SS_ROUND_ROBIN`: This scheduler provides a fair scheduling based on the number of user messages by cycling around non-empty stream queues.

`SCTP_SS_ROUND_ROBIN_PACKET`: This is a round-robin scheduler but only bundles user messages of the same stream in one packet. This minimizes head-of-line blocking when a packet is lost because only a single stream is affected.

`SCTP_SS_PRIORITY`: Scheduling with different priorities is used. Streams having a higher priority will be scheduled first and when multiple streams have the same priority, the default scheduling should be used for them. The priority can be assigned with the `sctp_stream_value` struct. The higher the assigned value, the lower the priority, that is the default value 0 is the highest priority and therefore the default scheduling will be used if no priorities have been assigned.

`SCTP_SS_FAIR_BANDWIDTH`: A fair bandwidth distribution between the streams can be activated using this value. This scheduler

considers the lengths of the messages of each stream and schedules them in a certain way to maintain an equal bandwidth for all streams.

SCTP_SS_FIRST_COME: The simple first-come, first-serve algorithm is selected by using this value. It just passes through the messages in the order in which they have been delivered by the application. No modification of the order is done at all.

4.1.3. Get or Set the Stream Scheduler Parameter (SCTP_SS_VALUE)

Some schedulers require additional information to be set for single streams as shown in the following table:

name	per stream info
SCTP_SS_DEFAULT	no
SCTP_SS_RR	no
SCTP_SS_RR_INTER	no
SCTP_SS_RR_PKT	no
SCTP_SS_RR_PKT_INTER	no
SCTP_SS_PRIO	yes
SCTP_SS_PRIO_INTER	yes
SCTP_SS_FB	no
SCTP_SS_FB_INTER	no
SCTP_SS_FCFS	no

This is achieved with the SCTP_SS_VALUE option and the corresponding struct sctp_stream_value. The definition of struct sctp_stream_value is as follows:

```
struct sctp_stream_value {
    sctp_assoc_t assoc_id;
    uint16_t stream_id;
    uint16_t stream_value;
};
```

assoc_id: Holds the identifier for the association of which the scheduler should be changed. The special SCTP_{FUTURE|CURRENT|ALL}_ASSOC can also be used. This parameter is ignored for one-to-one style sockets.

stream_id: Holds the stream id for the stream for which additional information has to be provided.

stream_value: The meaning of this field depends on the scheduler specified. It is ignored when the scheduler does not need additional information.

5. IANA Considerations

[NOTE to RFC-Editor:

"RFCXXXX" is to be replaced by the RFC number you assign this document.

]

[NOTE to RFC-Editor:

The suggested values for the chunk type and the chunk flags are tentative and to be confirmed by IANA.

]

This document (RFCXXXX) is the reference for all registrations described in this section.

A new chunk type has to be assigned by IANA. IANA should assign this value from the pool of chunks with the upper two bits set to '00'. This requires an additional line in the "Chunk Types" registry for SCTP:

ID Value	Chunk Type	Reference
17	New DATA chunk (N-DATA)	[RFCXXXX]

The registration table as defined in [RFC6096] for the chunk flags of this chunk type is initially given by the following table:

Chunk Flag Value	Chunk Flag Name	Reference
0x01	E bit	[RFCXXXX]
0x02	B bit	[RFCXXXX]
0x04	U bit	[RFCXXXX]
0x08	I bit	[RFCXXXX]

0x10	Unassigned		
0x20	Unassigned		
0x40	Unassigned		
0x80	Unassigned		
+-----+-----+-----+			

6. Security Considerations

This document does not add any additional security considerations in addition to the ones given in [RFC4960] and [RFC6458].

7. Acknowledgments

The authors wish to thank Lixia Zhang for her invaluable comments.

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.
- [RFC5061] Stewart, R., Xie, Q., Tuexen, M., Maruyama, S., and M. Kozuka, "Stream Control Transmission Protocol (SCTP) Dynamic Address Reconfiguration", RFC 5061, September 2007.
- [RFC6096] Tuexen, M. and R. Stewart, "Stream Control Transmission Protocol (SCTP) Chunk Flags Registration", RFC 6096, January 2011.
- [I-D.ietf-tsvwg-sctp-sack-immediately] Tuexen, M., Ruengeler, I., and R. Stewart, "SACK-IMMEDIATELY Extension for the Stream Control Transmission Protocol", draft-ietf-tsvwg-sctp-sack-immediately-04 (work in progress), August 2013.

8.2. Informative References

- [RFC3261] Rosenberg, J., Schulzrinne, H., Camarillo, G., Johnston, A., Peterson, J., Sparks, R., Handley, M., and E. Schooler, "SIP: Session Initiation Protocol", RFC 3261, June 2002.

[RFC6458] Stewart, R., Tuexen, M., Poon, K., Lei, P., and V.
Yasevich, "Sockets API Extensions for the Stream Control
Transmission Protocol (SCTP)", RFC 6458, December 2011.

Authors' Addresses

Randall R. Stewart
Adara Networks
Chapin, SC 29036
US

Email: randall@lakerest.net

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

Email: tuexen@fh-muenster.de

Salvatore Loreto
Ericsson
Hirsalantie 11
Jorvas 02420
FI

Email: Salvatore.Loreto@ericsson.com

Robin Seggelmann
T-Systems International GmbH
Fasanenweg 5
70771 Leinfelden-Echterdingen
DE

Email: robin.seggelmann@t-systems.com

Network Working Group
Internet-Draft
Intended status: Informational
Expires: April 23, 2014

M. Tuexen
Muenster Univ. of Appl. Sciences
R. Seggelmann
T-Systems International GmbH
R. Stewart
Adara Networks
S. Loreto
Ericsson
October 20, 2013

Additional Policies for the Partial Delivery Extension of the Stream
Control Transmission Protocol
draft-tuexen-tsvwg-sctp-prpolicies-03.txt

Abstract

This document defines policies for the Partial Reliability Extension of the Stream Control Transmission Protocol (PR-SCTP) allowing to limit the number of retransmissions or to prioritize user messages for more efficient send buffer usage.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on April 23, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents

carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
1.1. Overview	2
1.2. Data Types	2
2. Additional PR-SCTP Policies	3
2.1. Limited Retransmissions Policy	3
2.2. Priority Policy	3
3. Socket API Considerations	3
3.1. Support for Added PR-SCTP Policies	3
3.2. Socket Option for Getting the PR-SCTP Status (SCTP_PR_STATUS)	4
4. IANA Considerations	5
5. Security Considerations	5
6. Acknowledgments	5
7. References	5
7.1. Normative References	5
7.2. Informative References	5
Authors' Addresses	6

1. Introduction

1.1. Overview

The SCTP Partial Reliability Extension (PR-SCTP) defined in [RFC3758] provides a generic method for senders to abandon user messages. The decision to abandon a user message is sender side only and the exact condition is called a PR-SCTP policy. [RFC3758] also defines one particular PR-SCTP policy, called Timed Reliability. This allows the sender to specify a timeout for a user message after which the SCTP stack abandons the user message.

This document specifies two additional PR-SCTP policies:

Limited Retransmission Policy: Allows to limit the number of retransmissions.

Priority Policy: Allows to discard lower priority messages if space for higher priority messages is needed in the send buffer.

1.2. Data Types

This documents uses data types from Draft 6.6 (March 1997) of POSIX 1003.1g: `uintN_t` means an unsigned integer of exactly N bits (e.g. `uint16_t`). This is the same as in [RFC6458]

2. Additional PR-SCTP Policies

2.1. Limited Retransmissions Policy

Using the Limited Retransmission Policy allows the sender of a user message to specify an upper limit for the number of retransmissions for each DATA chunk of the given user messages. The sender must abandon a user message if the number of retransmissions of any of the DATA chunks of the user message would exceed the provided limit. Please note that the number of retransmissions includes the fast and the timer based retransmissions.

Limiting the number of retransmissions to 0 is allowed. This provides a service similar to UDP, which also does not send any retransmissions either.

The Limited Retransmissions Policy is used for data channels in the RTCWeb protocol stack.

2.2. Priority Policy

Using the Priority Policy allows the sender of a user message to specify a priority. When storing a user message in the send buffer while there is not enough available space, the SCTP stack may abandon other user messages with a priority lower than the provided one.

After lower priority messages have been abandoned high priority messages can be transferred without blocking the `send()` call.

The Priority Policy can be used in the IPFIX protocol stack. See [RFC7011] for more information.

3. Socket API Considerations

This section describes how the socket API defined in [RFC6458] is extended to support the newly defined PR-SCTP policies and to provide some statistical information.

Please note that this section is informational only.

3.1. Support for Added PR-SCTP Policies

As defined in [RFC6458], the PR-SCTP policy is specified and configured by using the following `sctp_prinfo` structure:

```

struct sctp_prinfo {
    uint16_t pr_policy;
    uint32_t pr_value;
};

```

When the Limited Retransmission Policy described in Section 2.1 is used, `pr_policy` has the value `SCTP_PR_SCTP_RTX` and the number of retransmissions is given in `pr_value`.

For using the Priority Policy described in Section 2.2, `pr_policy` has the value `SCTP_PR_SCTP_PRIO`. The priority is given in `pr_value`. The value of zero is the highest priority and larger numbers in `pr_value` denote lower priorities.

The following table summarizes the possible parameter settings defined in [RFC6458] and this document:

<code>pr_policy</code>	<code>pr_value</code>	Specification
<code>SCTP_PR_SCTP_NONE</code>	Ignored	[RFC6458]
<code>SCTP_PR_SCTP_TTL</code>	Lifetime in ms	[RFC6458]
<code>SCTP_PR_SCTP_RTX</code>	Number of retransmissions	Section 2.1
<code>SCTP_PR_SCTP_PRIO</code>	Priority	Section 2.2

3.2. Socket Option for Getting the PR-SCTP Status (`SCTP_PR_STATUS`)

This socket option uses `IPPROTO_SCTP` as its level and `SCTP_PR_STATUS` as its name. It can only be used with `getsockopt()`, but not with `setsockopt()`. The socket option value uses the following structure:

```

struct sctp_prstatus {
    sctp_assoc_t sprstat_assoc_id;
    uint32_t sprstat_abandoned_unsent;
    uint32_t sprstat_abandoned_sent;
};

```

`sprstat_assoc_id`: This parameter is ignored for one-to-one style sockets. For one-to-many style sockets this parameter indicates for which association the user wants the information. It is an error to use `SCTP_{CURRENT|ALL|FUTURE}_ASSOC` in `sprstat_assoc_id`

`sprstat_abandoned_unsent`: The number of user messages which have been abandoned, before any part of the user message could be sent.

`sprstat_abandoned_sent`: The number of user messages which have been abandoned, after a part of the user message has been sent.

There are separate counters for unsent and sent user messages because the `SCTP_SEND_FAILED_EVENT` supports a similar differentiation. Please note that an abandoned large user messages requiring an SCTP level fragmentation is reported in the `sprstat_abandoned_sent` counter as soon as at least one fragment of it has been sent. Therefore each abandoned user messages is either counted in `sprstat_abandoned_unsent` or `sprstat_abandoned_sent`.

If more detailed information about abandoned user messages is required, the subscription to the `SCTP_SEND_FAILED_EVENT` is recommended.

`sctp_opt_info()` needs to be extended to support `SCTP_PR_STATUS`.

4. IANA Considerations

This document requires no actions from IANA.

5. Security Considerations

This document does not add any additional security considerations in addition to the ones given in [RFC4960], [RFC3758], and [RFC6458].

6. Acknowledgments

The authors wish to thank Irene Ruengeler, Jamal Hadi Salim, and Vlad Yasevich for there invaluable comments.

7. References

7.1. Normative References

- [RFC3758] Stewart, R., Ramalho, M., Xie, Q., Tuexen, M., and P. Conrad, "Stream Control Transmission Protocol (SCTP) Partial Reliability Extension", RFC 3758, May 2004.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", RFC 4960, September 2007.

7.2. Informative References

- [RFC6458] Stewart, R., Tuexen, M., Poon, K., Lei, P., and V. Yasevich, "Sockets API Extensions for the Stream Control Transmission Protocol (SCTP)", RFC 6458, December 2011.

[RFC7011] Claise, B., Trammell, B., and P. Aitken, "Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information", STD 77, RFC 7011, September 2013.

Authors' Addresses

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

Email: tuexen@fh-muenster.de

Robin Seggelmann
T-Systems International GmbH
Fasanenweg 5
70771 Leinfelden-Echterdingen
DE

Email: robin.seggelmann@t-systems.com

Randall R. Stewart
Adara Networks
Chapin, SC 29036
US

Email: randall@lakerest.net

Salvatore Loreto
Ericsson
Hirsalantie 11
Jorvas 02420
FI

Email: Salvatore.Loreto@ericsson.com

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: April 24, 2014

E. Crabbe, Ed.
Google
L. Yong, Ed.
Huawei USA
X. Xu, Ed.
Huawei Technologies
October 21, 2013

Generic UDP Encapsulation for IP Tunneling
draft-yong-tsvwg-gre-in-udp-encap-02

Abstract

This document describes a method of encapsulating arbitrary protocols within GRE and UDP headers. In this encapsulation, the source UDP port may be used as an entropy field for purposes of loadbalancing while the payload protocol may be identified by the GRE Protocol Type.

Requirements Language

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].

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on April 24, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	2
2. Terminology	3
3. Procedures	3
4. Encapsulation Considerations	6
5. Backward Compatibility	7
6. IANA Considerations	7
7. Security Considerations	7
7.1. Vulnerability	7
8. Acknowledgements	8
9. Contributing Authors	8
10. References	9
10.1. Normative References	9
10.2. Informative References	10
Authors' Addresses	10

1. Introduction

Load balancing, or more specifically, statistical multiplexing of traffic using Equal Cost Multi-Path (ECMP) and/or Link Aggregation Groups (LAGs) in IP networks is a widely used technique for creating higher capacity networks out of lower capacity links. Most existing routers in IP networks are already capable of distributing IP traffic flows over ECMP paths and/or LAGs on the basis of a hash function performed on flow invariant fields in IP packet headers and their payload protocol headers. Specifically, when the IP payload is a User Datagram Protocol (UDP)[RFC0768] or Transmission Control Protocol (TCP) packet, router hash functions frequently operate on the five-tuple of the source IP address, the destination IP address, the source port, the destination port, and the protocol/next-header

Several tunneling techniques are in common use in IP networks, such as Generic Routing Encapsulation (GRE) [RFC2784], MPLS [RFC4023] and L2TPv3 [RFC3931]. GRE is an increasingly popular encapsulation choice, especially in environments where MPLS is unavailable or unnecessary. Unfortunately, use of common GRE endpoints may reduce the entropy available for use in load balancing, especially in environments where the GRE Key field [RFC2890] is not readily available for use as entropy in forwarding decisions.

This document defines a generic GRE-in-UDP encapsulation for tunneling arbitrary network protocol payloads across an IP network environment where ECMP or LAGs are used. The GRE header provides payload protocol de-multiplexing by way of it's protocol type field [RFC2784] while the UDP header provides additional entropy by way of it's source port.

This encapsulation method requires no changes to the transit IP network. Hash functions in most existing IP routers may utilize and benefit from the use of a GRE-in-UDP tunnel without needing any change or upgrade to their ECMP implementations. The encapsulation mechanism is applicable to a variety of IP networks including Data Center and wide area networks.

2. Terminology

The terms defined in [RFC0768] are used in this document.

3. Procedures

When a tunnel ingress device conforming to this document receives a packet, the ingress MUST encapsulate the packet in UDP and GRE headers and set the destination port of the UDP header to [TBD] Section 6. The ingress device must also insert the payload protocol type in the GRE Protocol Type field. The ingress device SHOULD set the UDP source port based on flow invariant fields from the payload header, otherwise it should be set to a randomly selected constant value, e.g. zero, to avoid packet flow reordering. How a tunnel ingress generates entropy from the payload is outside the scope of this document. The tunnel ingress MUST encode its own IP address as the source IP address and the egress tunnel endpoint IP address. The TTL field in the IP header must be set to a value appropriate for delivery of the encapsulated packet to the tunnel egress endpoint.

When the tunnel egress receives a packet, it must remove the outer UDP and GRE headers. Section 5 describes the error handling when this entity is not instantiated at the tunnel egress.

To simplify packet processing at the tunnel egress, packets destined to this assigned UDP destination port [TBD] SHOULD have their UDP checksum and Sequence flags set to zero because the egress tunnel only needs to identify this protocol. Although IPv6 [RFC2460] restricts the processing a packet with the UDP checksum of zero, [RFC6935] and [RFC6936] relax this constraint to allow the zero UDP checksum.

The tunnel ingress may set the GRE Key Present, Sequence Number Present, and Checksum Present bits and associated fields in the GRE header defined by [RFC2784] and [RFC2890].

In addition IPv6 nodes MUST conform to the following:

1. the IPv6 tunnel ingress and egress SHOULD follow the node requirements specified in Section 4 of [RFC6936] and the usage requirements specified in Section 5 of [RFC6936]
2. IPv6 transit nodes SHOULD follow the requirements 9, 10, 11 specified in Section 5 of [RFC6936].

The format of the GRE-in-UDP encapsulation for both IPv4 and IPv6 outer headers is shown in the following figures:

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

```

IPv4 Header:

```

+-----+-----+-----+-----+-----+-----+-----+-----+
|Version|  IHL  |Type of Service|          Total Length          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          Identification          |Flags|      Fragment Offset      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  Time to Live |Protocol=17[UDP]|          Header Checksum          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Source IPv4 Address                               |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Destination IPv4 Address                               |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

UDP Header:

```

+-----+-----+-----+-----+-----+-----+-----+-----+
|          Source Port = XXXX          |          Dest Port = TBD          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          UDP Length          |          UDP Checksum          |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

GRE Header:

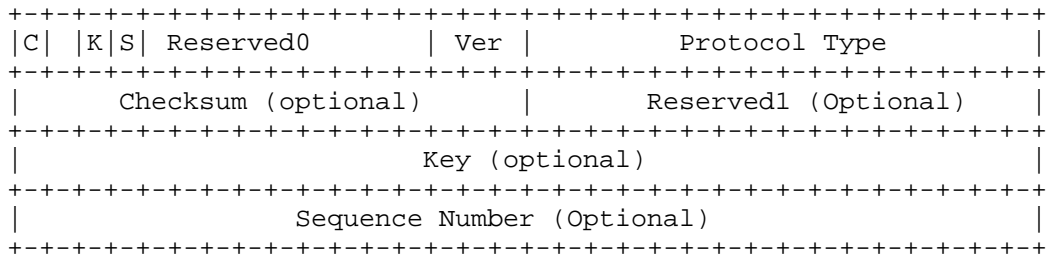
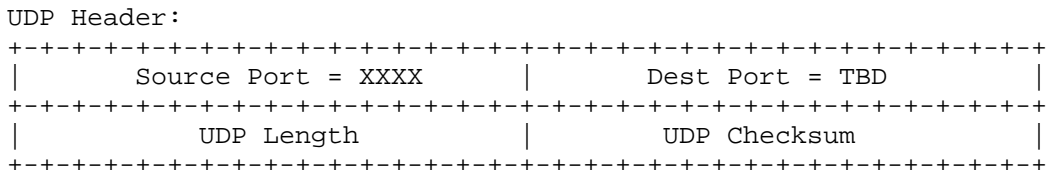
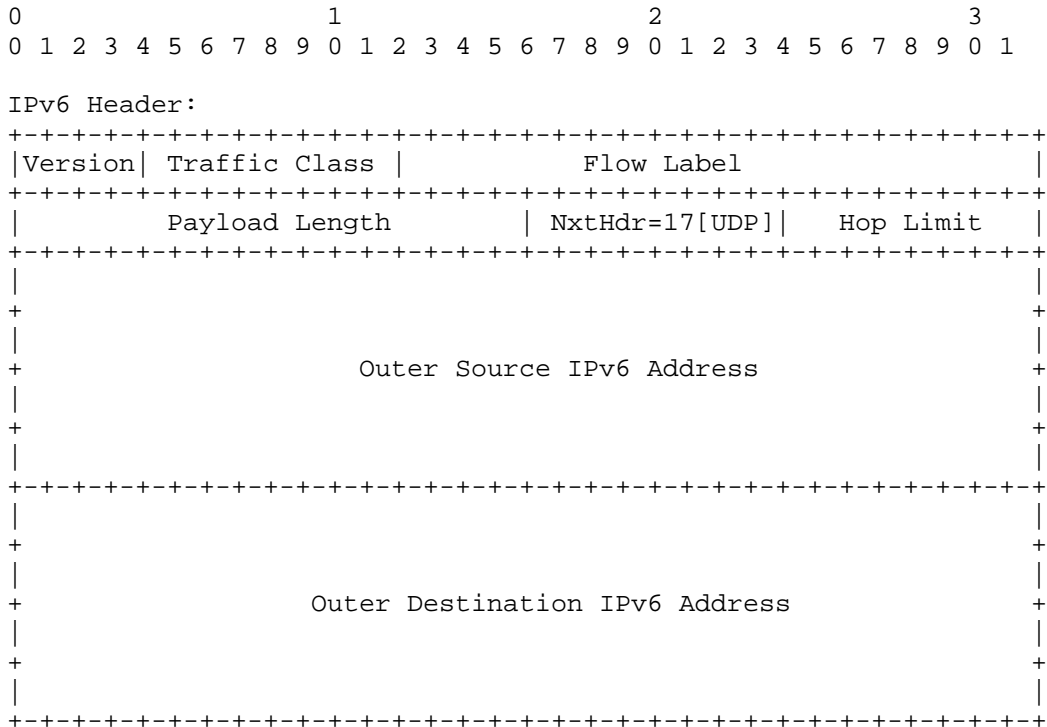


Figure 1: UDP+GRE IPv4 headers



GRE Header:

C	K S	Reserved0		Ver		Protocol Type	
Checksum (optional)				Reserved1 (Optional)			
Key (optional)							
Sequence Number (Optional)							

Figure 2: UDP+GRE IPv6 headers

The total overhead increase for a UDP+GRE tunnel without use of optional GRE fields, representing the lowest total overhead increase, is 32 bytes in the case of IPv4 and 52 bytes in the case of IPv6. The total overhead increase for a UDP+GRE tunnel with use of GRE Key, Sequence and Checksum Fields, representing the highest total overhead increase, is 44 bytes in the case of IPv4 and 64 bytes in the case of IPv6.

4. Encapsulation Considerations

GRE-in-UDP encapsulation allows the tunneled traffic to be unicast, broadcast, or multicast traffic. Entropy may be generated from the header of tunneled unicast or broadcast/multicast packets at tunnel ingress. The mapping mechanism between the tunneled multicast traffic and the multicast capability in the IP network is transparent and independent to the encapsulation and is outside the scope of this document.

If tunnel ingress must perform fragmentation on a packet before encapsulation, it MUST use the same source UDP port for all packet fragments. This ensures that the transit routers will forward the packet fragments on the same path. GRE-in-UDP encapsulation introduces some overhead as mentioned in section 3, which reduces the effective Maximum Transmission Unit (MTU) size. An operator should factor in this addition overhead bytes when considering an MTU size for the payload to reduce the likelihood of fragmentation.

To ensure the tunneled traffic gets the same treatment over the IP network, prior to the encapsulation process, tunnel ingress should process the payload to get the proper parameters to fill into the IP header such as DiffServ [[RFC2983]]. Tunnel end points that support ECN MUST use the method described in [RFC6040] for ECN marking propagation. This process is outside of the scope of this document.

Note that the IPv6 header [RFC2460] contains a flow label field that may be used for load balancing in an IPv6 network [RFC6438]. Thus in an IPv6 network, either GRE-in-UDP or flow labels may be used in order to improve load balancing performance. Use of GRE-in-UDP encapsulation provides a unified hardware implementation for load balancing in an IP network independent of the IP version(s) in use.

5. Backward Compatibility

It is assumed that tunnel ingress routers must be upgraded in order to support the encapsulations described in this document.

No change is required at transit routers to support forwarding of the encapsulation described in this document.

If a router that is intended for use as a tunnel egress does not support the GRE-in-UDP encapsulation described in this document, it will not be listening on destination port [TBD]. In these cases, the router will conform to normal UDP processing and respond to the tunnel ingress with an ICMP message indicating "port unreachable" according to [RFC0792]. Upon receiving this ICMP message, the tunnel ingress MUST NOT continue to use GRE-in-UDP encapsulation toward this tunnel egress without management intervention.

6. IANA Considerations

IANA is requested to make the following allocation: Service Name: GRE-in-UDP Transport Protocol(s): UDP Assignee: IESG iesg@ietf.org Contact: IETF Chair chair@ietf.org Description: GRE-in-UDP Encapsulation Reference: [This.I-D] Port Number: TBD Service Code: N/A A Known Unauthorized Uses: N/A Assignment Notes: N/A

7. Security Considerations

7.1. Vulnerability

Neither UDP nor GRE encapsulation effects security for the payload protocol. When using GRE-in-UDP, Network Security in a network is similar to that of a network using GRE.

Use of ICMP for signaling of the GRE-in-UDP encapsulation capability adds a security concern. Tunnel ingress devices may want to validate the origin of ICMP Port Unreachable messages before taking action. The mechanism for performing this validation is out of the scope of this document.

In an instance where the UDP src port is not set based on the flow invariant fields from the payload header, a random port SHOULD be

selected in order to minimize the vulnerability to off-path attacks.
[RFC6056] How the src port randomization occurs is outside scope of
this document.

8. Acknowledgements

The Authors would like to thank Vivek Kumar, Ron Bonica, Joe Touch,
Ruediger Geib, Gorrry Fairhurst, and David Black for their review and
valuable input on this draft.

9. Contributing Authors

The following people all contributed significantly to this document
and are listed below in alphabetical order:

John E. Drake
Juniper Networks

Email: jdrake@juniper.net

Adrian Farrel
Juniper Networks

Email: adrian@olddog.co.uk

Vishwas Manral
Hewlett-Packard Corp.
3000 Hanover St, Palo Alto.

Email: vishwas.manral@hp.com

Carlos Pignataro
Cisco Systems
7200-12 Kit Creek Road
Research Triangle Park, NC 27709 USA

EMail: cpignata@cisco.com

Yongbing Fan
China Telecom
Guangzhou, China.
Phone: +86 20 38639121

10. References

10.1. Normative References

- [RFC0768] Postel, J., "User Datagram Protocol", STD 6, RFC 768, August 1980.
- [RFC0791] Postel, J., "Internet Protocol", STD 5, RFC 791, September 1981.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2784] Farinacci, D., Li, T., Hanks, S., Meyer, D., and P. Traina, "Generic Routing Encapsulation (GRE)", RFC 2784, March 2000.
- [RFC2890] Dommety, G., "Key and Sequence Number Extensions to GRE", RFC 2890, September 2000.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", RFC 2983, October 2000.
- [RFC5405] Eggert, L. and G. Fairhurst, "Unicast UDP Usage Guidelines for Application Designers", BCP 145, RFC 5405, November 2008.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, November 2010.
- [RFC6335] Cotton, M., Eggert, L., Touch, J., Westerlund, M., and S. Cheshire, "Internet Assigned Numbers Authority (IANA) Procedures for the Management of the Service Name and Transport Protocol Port Number Registry", BCP 165, RFC 6335, August 2011.
- [RFC6438] Carpenter, B. and S. Amante, "Using the IPv6 Flow Label for Equal Cost Multipath Routing and Link Aggregation in Tunnels", RFC 6438, November 2011.
- [RFC6935] Eubanks, M., Chimento, P., and M. Westerlund, "IPv6 and UDP Checksums for Tunneled Packets", RFC 6935, April 2013.
- [RFC6936] Fairhurst, G. and M. Westerlund, "Applicability Statement for the Use of IPv6 UDP Datagrams with Zero Checksums", RFC 6936, April 2013.

10.2. Informative References

- [RFC0792] Postel, J., "Internet Control Message Protocol", STD 5, RFC 792, September 1981.
- [RFC2460] Deering, S. and R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification", RFC 2460, December 1998.
- [RFC3931] Lau, J., Townsley, M., and I. Goyret, "Layer Two Tunneling Protocol - Version 3 (L2TPv3)", RFC 3931, March 2005.
- [RFC4023] Worster, T., Rekhter, Y., and E. Rosen, "Encapsulating MPLS in IP or Generic Routing Encapsulation (GRE)", RFC 4023, March 2005.
- [RFC4364] Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)", RFC 4364, February 2006.
- [RFC4884] Bonica, R., Gan, D., Tappan, D., and C. Pignataro, "Extended ICMP to Support Multi-Part Messages", RFC 4884, April 2007.
- [RFC6790] Kompella, K., Drake, J., Amante, S., Henderickx, W., and L. Yong, "The Use of Entropy Labels in MPLS Forwarding", RFC 6790, November 2012.

Authors' Addresses

Edward Crabbe (editor)
Google
1600 Amphitheatre Parkway
Mountain View, CA 94102
US

Email: edward.crabbe@gmail.com

Lucy Yong (editor)
Huawei USA
5340 Legacy Drive
San Jose, TX 75025
US

Email: lucy.yong@huawei.com

Xiaohu Xu (editor)
Huawei Technologies
Beijing
China

Email: xuxiaohu@huawei.com

Internet Engineering Task Force
Internet-Draft
Intended status: Standards Track
Expires: January 10, 2014

T. Eckert, Ed.
A. Zamfir
A. Choukir
Cisco
July 09, 2013

Flow Metadata Signaling with RSVP
draft-zamfir-tsvwg-flow-metadata-rsvp-00

Abstract

This specification proposes RSVP protocol extensions for signaling flow metadata attributes.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on January 10, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	2
1.1. Requirements Language	2
2. Flow Metadata Object	2
2.1. FLOW_METADATA Class	3
2.1.1. Basic IPFIX FLOW_METADATA Object	3
2.1.2. Enhanced Protocol Independent FLOW_METADATA Object	3
2.2. Semantic of carrying the Metadata Object	3
2.3. Processing by a Non-Metadata Capable RSVP Router	4
2.4. Processing by a Metadata Capable RSVP Router	4
3. References	5
3.1. Normative References	5
3.2. Informative References	5
Authors' Addresses	6

1. Introduction

Flow Metadata attributes are information elements (attributes) that identify flow characteristics, such as the type of media carried by application flows (e.g. video), the service class, the application that originated the flow, and others. The description of the Flow Metadata technology and some of the attribute definitions can be found in [I-D.eckert-intarea-flow-metadata-framework]. The flow attributes can be signaled over the flow path and inspected by intermediate network nodes for the purpose of applying differentiated flow treatment or collect network analytics. This specification proposes the use of RSVP as signaling protocol to carry the Flow Metadata using a new RSVP object. Two C-Type values are proposed for this object to allow for two possible encodings.

1.1. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC2119].

2. Flow Metadata Object

This specification proposes a new RSVP object with Class-Num from the 0x1b bbbbbb range. To support informational metadata attribute processing on the path to the receiver, the sender inserts the Metadata object into an IPv4 or IPv6 Path message (i.e. Path messages with SESSION Class = 1 and SENDER_TEMPLATE Class = 11). The Metadata object SHOULD appear only once in the message.

The object definition is given in Section 2.1 while the details of processing are covered in Section 2.2

2.1. FLOW_METADATA Class

FLOW_METADATA Class = 234

Two encodings are defined, both of which carry the same IPFIX registered attributes as defined in [I-D.eckert-intarea-flow-metadata-framework]. The first encoding (Basic IPFIX FLOW_METADATA) has less flexibility and lower encoding efficiency. This version of the encoding is referenced here for legacy reasons. It does not support a range of options that the second one does, including the signaling of sender and receiver attributes, security elements, distinction of originator of the attributes and ease of extensibility.

2.1.1. Basic IPFIX FLOW_METADATA Object

Basic IPFIX FLOW_METADATA Object: Class = 234, C-Type = 1

- o The metadata attributes are encoded in IPFIX format, as described in [RFC5101], with the following restrictions when creating the object:
 - * Options Template Record MUST NOT be present
 - * One and only one Template Record MUST be present
 - * One and only one Data Record MUST be included
- o An intermediate node that supports this specification SHOULD ignore any Options Template Record. It SHOULD only decode and process the first occurring Template and Data Records.

2.1.2. Enhanced Protocol Independent FLOW_METADATA Object

Enhanced Protocol Independent FLOW_METADATA Object: Class = 234, C-Type = 2

- o The contents and encoding rules for this object are specified in [I-D.eckert-intarea-flow-metadata-framework] and [I-D.choukir-tsvwg-flow-metadata-encoding].

2.2. Semantic of carrying the Metadata Object

The Metadata Object included in the Path message carries attributes from the sender of the flow towards the receiver. In some cases, e.g. if the sender does not support the generation and signaling of Metadata attribute, these attributes may be inserted by a proxy along the path of the flow. Metadata RSVP nodes on path may modify the

metadata attributes for purpose of influencing policy toward the receiver.

The node that originates Metadata information in a Path message may do so for the sole purpose of signaling Metadata information. In this case, the SENDER_TSPEC objects fields (as defined by [RFC2210]) should be set to 0:

- o Token Bucket Rate [r]
- o Token Bucket Size [b]
- o Peak Data Rate [p]
- o Minimum Policed Unit [m]

If the Metadata object is inserted in a Path message used for IntServ service [[RFC2210]] reservation requests, then all the rules of RSVP reservation request apply and in addition any actions driven purely by the metadata attributes may equally take place.

While the Metadata Object may be included in a Resv message, the specific processing rules for this option is left for followup documents or future versions of this specification.

2.3. Processing by a Non-Metadata Capable RSVP Router

As described in [RFC2205], a node that does not understand the Metadata object, should ignore but forward it, unexamined and unmodified. When received in Path or Resv messages, it should be saved with the corresponding state and forwarded in any refresh message resulting from that state.

2.4. Processing by a Metadata Capable RSVP Router

The Metadata object may be inserted by the data flow initiating endpoint or network nodes along the path. The means by which an implementation determines the content of the Metadata object is outside the scope of this document.

Intermediate nodes that support this specification, decode the Flow Metadata information as indicated by the C-Type field only when received in Path message. Depending on the attributes, local configuration and policies, the node may take some actions. The Metadata attribute semantics are described in [I-D.eckert-intarea-flow-metadata-framework]. The received Flow Metadata object is stored against the Path state. When a subsequent Path message is received with a modified Metadata object, the

intermediate node determines the attributes that have been removed, modified and/or added by comparing the old and new objects, and takes appropriate actions.

As a result of these actions, an intermediate node may add new attributes to the Metadata object received in the Path message and signal them downstream. It can also modify some of the attributes present in the Flow Metadata object. RSVP does not have any transport protocol specific restrictions and the exact set of attributes that can be inserted and modified by intermediate nodes is described in [I-D.eckert-intarea-flow-metadata-framework]. Depending on local policies, an intermediate node may also remove some of the attributes received in the Metadata object of a Path message before forwarding downstream.

An intermediate node that receives a Resv message with a Metadata Object SHOULD store the object against the state and forward it unexamined and unmodified.

3. References

3.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2205] Braden, B., Zhang, L., Berson, S., Herzog, S., and S. Jamin, "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification", RFC 2205, September 1997.
- [RFC2210] Wroclawski, J., "The Use of RSVP with IETF Integrated Services", RFC 2210, September 1997.
- [RFC5101] Claise, B., "Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of IP Traffic Flow Information", RFC 5101, January 2008.

3.2. Informative References

- [I-D.choukir-tsvwg-flow-metadata-encoding]
Eckert, T., Zamfir, A., Choukir, A., and C. Eckel,
"Protocol Independent Encoding for Signaling Flow
Characteristics", draft-choukir-tsvwg-flow-metadata-
encoding-00 (work in progress), July 2013.
- [I-D.eckert-intarea-flow-metadata-framework]
Eckert, T., Penno, R., Choukir, A., and C. Eckel, "A
Framework for Signaling Flow Characteristics between

Applications and the Network", draft-eckert-intarea-flow-
metadata-framework-00 (work in progress), July 2013.

Authors' Addresses

Toerless Eckert (editor)
Cisco Systems, Inc.
San Jose
US

Email: eckert@cisco.com

Anca Zamfir
Cisco Systems, Inc.
EPFL, Quartier de l'Innovation
Ecublens, Vaud 1015
Switzerland

Email: ancaz@cisco.com

Amine Choukir
Cisco Systems, Inc.
EPFL, Quartier de l'Innovation
Ecublens, Vaud 1015
CH

Phone: +41 78 75 98 561
Email: amchouki@cisco.com