

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

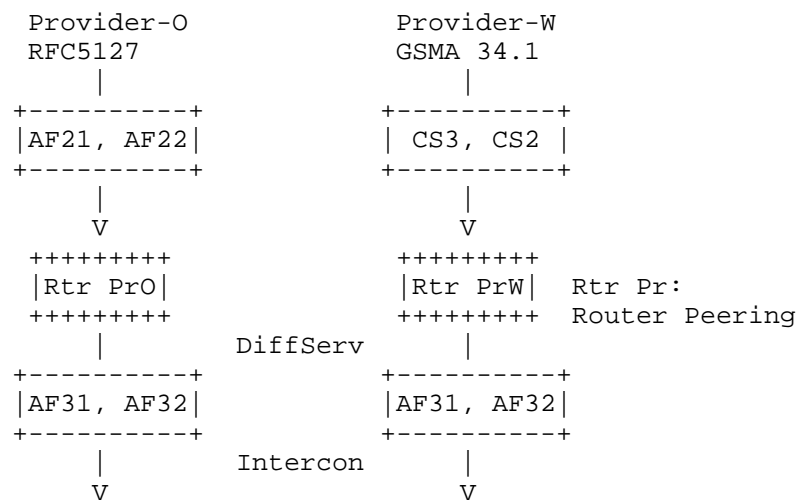
Network Control PHB treatment see section 3.2). GSMA IR.34 proposes to transport signaling traffic by AF31 too.

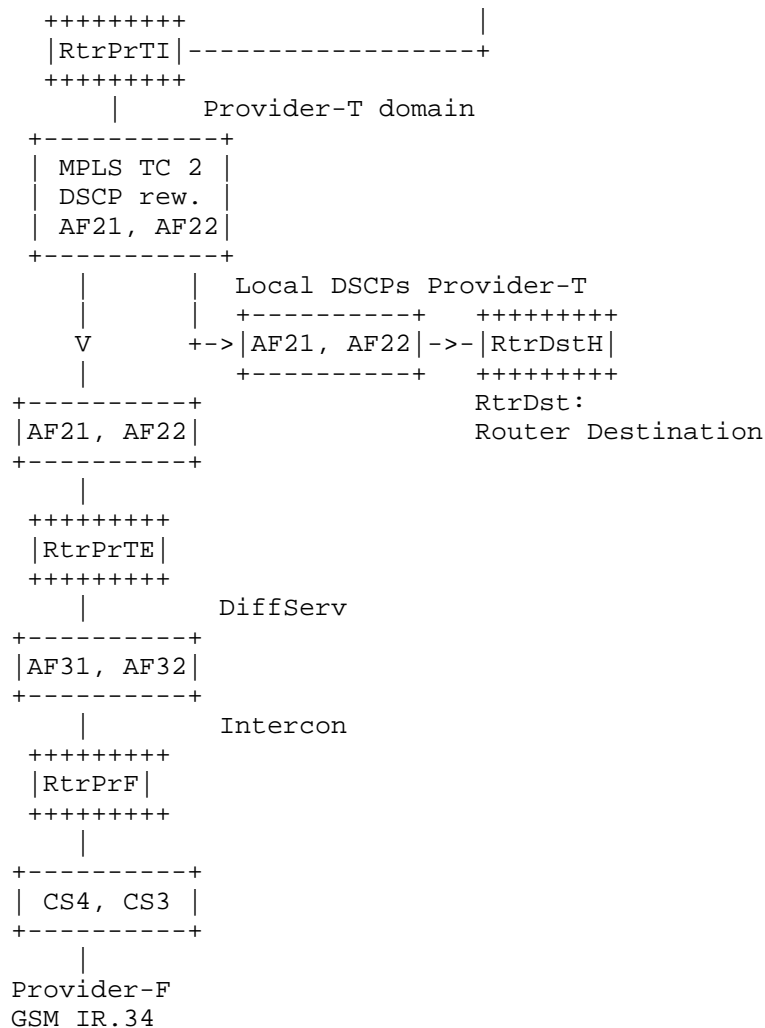
Default / Elastic Treatment Aggregate: transports the default PHB, CS0 with DSCP 000 000. RFC 5127 example refers to this Treatment Aggregate as Aggregate Elastic. An important difference as compared to RFC5127 is that any traffic with unrecognized or unsupported DSCPs may be remarked to this DSCP.

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 then 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 PHB's 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

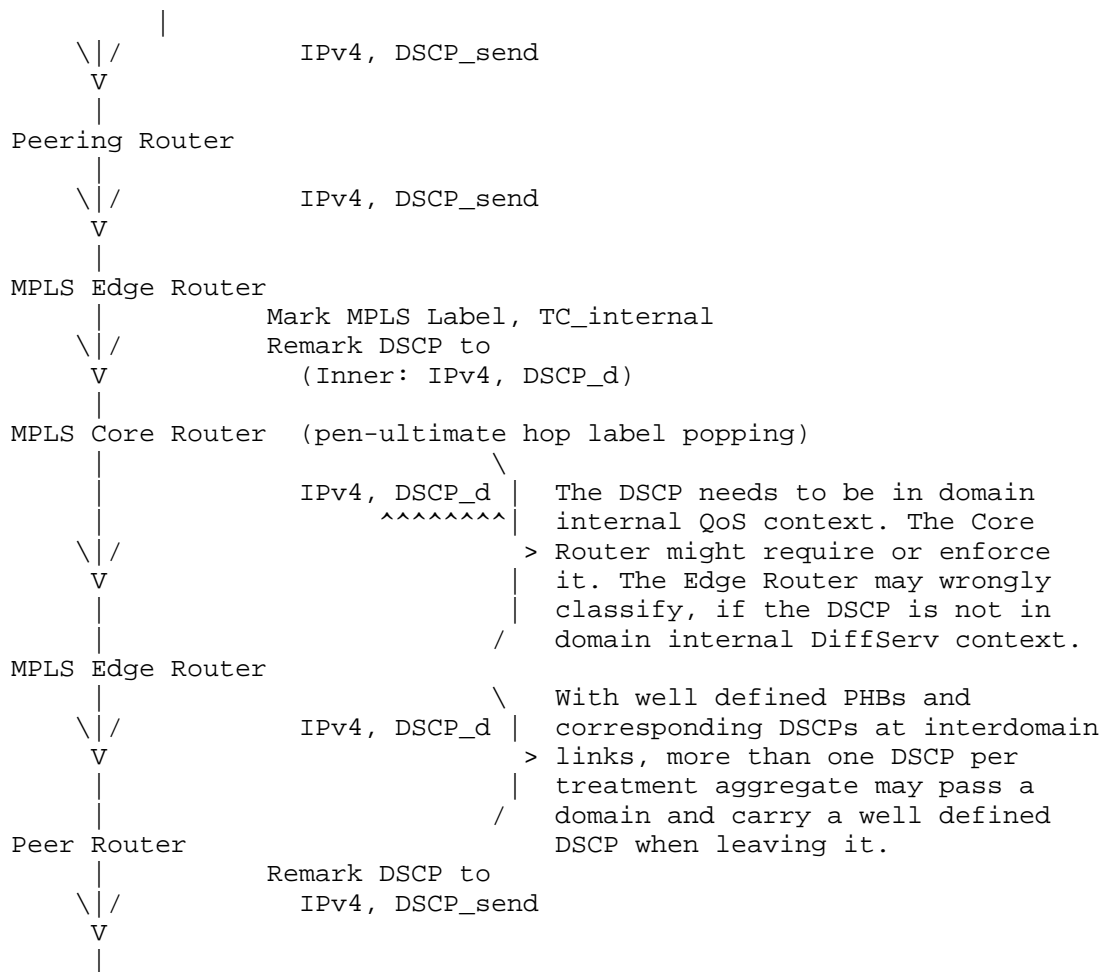
The MPLS Short Pipe Model (or Pen-ultimate Hop Label Popping) is widely deployed by IP carriers. If non-tunneled IPv4 traffic is transported using MPLS Short Pipe, IP headers appear inside the last section of the MPLS domain. This likely impacts the number of PHBs and DSCPs a network provider supports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

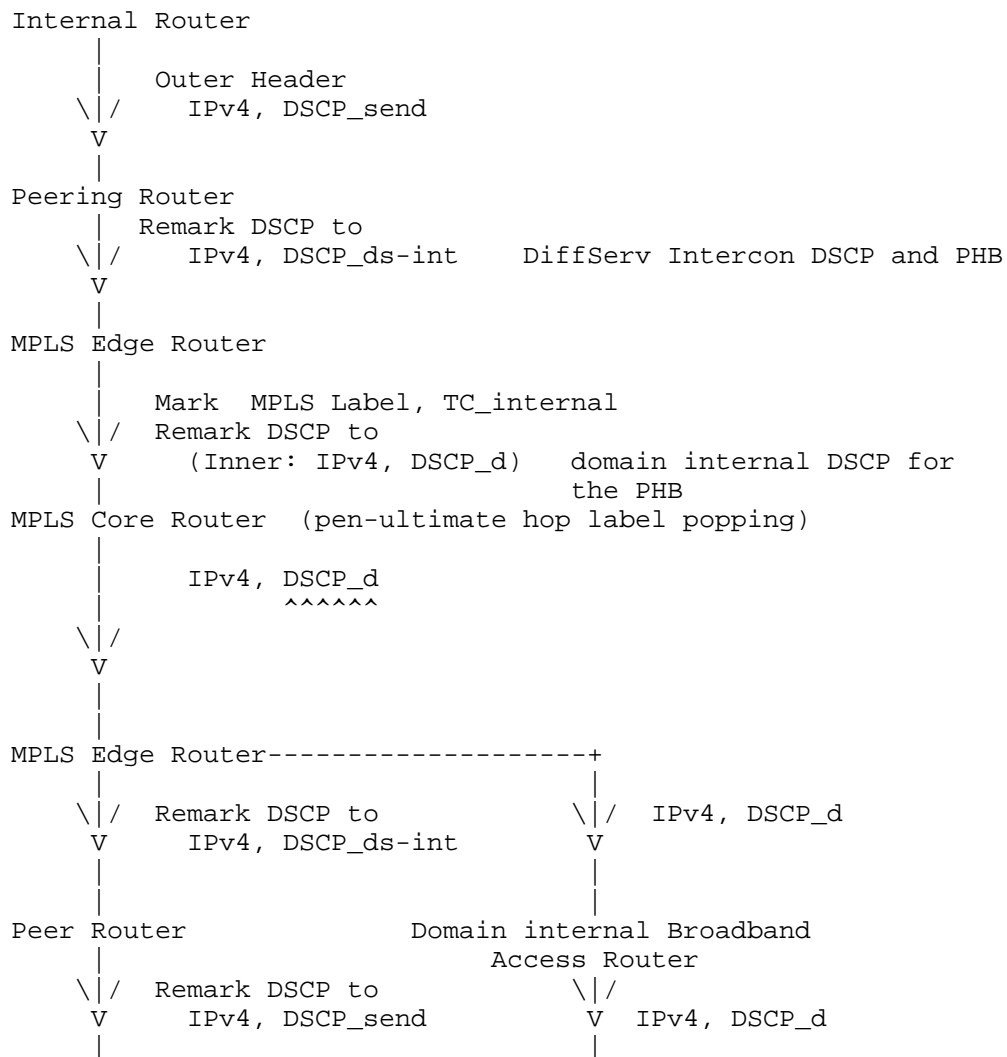
Figure 2

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