

Internet Engineering Task Force (IETF)
Request for Comments: 6660
Obsoletes: 5696
Category: Standards Track
ISSN: 2070-1721

B. Briscoe
BT
T. Moncaster
University of Cambridge
M. Menth
University of Tuebingen
July 2012

Encoding Three Pre-Congestion Notification (PCN) States
in the IP Header Using a Single Diffserv Codepoint (DSCP)

Abstract

The objective of Pre-Congestion Notification (PCN) is to protect the quality of service (QoS) of inelastic flows within a Diffserv domain. The overall rate of PCN-traffic is metered on every link in the PCN-domain, and PCN-packets are appropriately marked when certain configured rates are exceeded. Egress nodes pass information about these PCN-marks to Decision Points that then decide whether to admit or block new flow requests or to terminate some already admitted flows during serious pre-congestion.

This document specifies how PCN-marks are to be encoded into the IP header by reusing the Explicit Congestion Notification (ECN) codepoints within a PCN-domain. The PCN wire protocol for non-IP protocol headers will need to be defined elsewhere. Nonetheless, this document clarifies the PCN encoding for MPLS in an informational appendix. The encoding for IP provides for up to three different PCN marking states using a single Diffserv codepoint (DSCP): not-marked (NM), threshold-marked (ThM), and excess-traffic-marked (ETM). Hence, it is called the 3-in-1 PCN encoding. This document obsoletes RFC 5696.

Status of This Memo

This is an Internet Standards Track document.

This document is a product of the Internet Engineering Task Force (IETF). It represents the consensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Further information on Internet Standards is available in Section 2 of RFC 5741.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at <http://www.rfc-editor.org/info/rfc6660>.

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 Language	4
2.	Definitions and Abbreviations	4
2.1.	Terminology	4
2.2.	List of Abbreviations	5
3.	Definition of 3-in-1 PCN Encoding	6
4.	Requirements for and Applicability of 3-in-1 PCN Encoding	7
4.1.	PCN Requirements	7
4.2.	Requirements Imposed by Tunnelling	7
4.3.	Applicable Environments for the 3-in-1 PCN Encoding	8
5.	Behaviour of a PCN-node to Comply with the 3-in-1 PCN Encoding	8
5.1.	PCN-Ingress-Node Behaviour	8
5.2.	PCN-Interior-Node Behaviour	11
5.2.1.	Behaviour Common to All PCN-Interior-Nodes	11
5.2.2.	Behaviour of PCN-Interior-Nodes Using Two PCN-Markings	11
5.2.3.	Behaviour of PCN-Interior-Nodes Using One PCN-Marking	12
5.3.	PCN-Egress-Node Behaviour	13
6.	Backward Compatibility	13
6.1.	Backward Compatibility with ECN	13
6.2.	Backward Compatibility with the Encoding in RFC 5696	14
7.	Security Considerations	14
8.	Conclusions	15
9.	Acknowledgements	15
10.	References	15
10.1.	Normative References	15
10.2.	Informative References	16
	Appendix A. Choice of Suitable DSCPs	18
	Appendix B. Coexistence of ECN and PCN	19

Appendix C. Example Mapping between Encoding of PCN-Marks in IP and in MPLS Shim Headers	22
Appendix D. Rationale for Difference between the Schemes Using One PCN-Marking	23

1. Introduction

The objective of Pre-Congestion Notification (PCN) [RFC5559] is to protect 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, to decide whether to admit or block a new flow request, and flow termination to terminate some existing flows during serious pre-congestion. To achieve this, 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 real congestion occurs (hence "pre-congestion notification").

[RFC5670] provides for two metering and marking functions that are generally configured with different reference rates. Threshold-marking marks all PCN packets once their traffic rate on a link exceeds the configured reference rate (PCN-threshold-rate). Excess-traffic-marking marks only those PCN packets that exceed the configured reference rate (PCN-excess-rate). The PCN-excess-rate is typically larger than the PCN-threshold-rate [RFC5559]. Egress nodes monitor the PCN-marks of received PCN-packets and pass information about these PCN-marks to Decision Points that then decide whether to admit new flows or terminate existing flows [RFC6661] [RFC6662].

The encoding defined in [RFC5696] described how two PCN marking states (not-marked and PCN-marked) could be encoded into the IP header using a single Diffserv codepoint. It defined '01' as an experimental codepoint (EXP), along with guidelines for its use. Two PCN marking states are sufficient for the Single Marking edge behaviour [RFC6662]. However, PCN-domains utilising the controlled load edge behaviour [RFC6661] require three PCN marking states. This document extends the encoding that originally appeared in RFC 5696 by redefining the experimental codepoint as a third PCN marking state in the IP header, but still using a single Diffserv codepoint. This encoding scheme is therefore called the "3-in-1 PCN encoding". It obsoletes the [RFC5696] encoding, which provides only a subset of the same capabilities.

The full version of the 3-in-1 encoding requires any tunnel endpoint within the PCN-domain to support the normal tunnelling rules defined in [RFC6040]. There is one limited exception to this constraint

where the PCN-domain only uses the excess-traffic-marking behaviour and where the threshold-marking behaviour is deactivated. This is discussed in Section 5.2.3.1.

This document only concerns the PCN wire protocol encoding for IP headers, whether IPv4 or IPv6. It makes no changes or recommendations concerning algorithms for congestion marking or congestion response. Other documents will define the PCN wire protocol for other header types. Appendix C discusses a possible mapping between IP and MPLS.

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. Definitions and Abbreviations

2.1. Terminology

The terms PCN-domain, PCN-node, PCN-interior-node, PCN-ingress-node, PCN-egress-node, PCN-boundary-node, PCN-traffic, PCN-packets, and PCN-marking are used as defined in [RFC5559]. The following additional terms are defined in this document:

PCN encoding: mapping of PCN marking states to specific codepoints in the packet header.

PCN-compatible Diffserv codepoint: a Diffserv codepoint indicating packets for which the ECN field carries PCN-markings rather than [RFC3168] markings. Note that an operator configures PCN-nodes to recognise PCN-compatible DSCPs, whereas the same DSCP has no PCN-specific meaning to a node outside the PCN-domain.

Threshold-marked codepoint: a codepoint that indicates a packet has been threshold-marked; that is, a packet that has been marked at a PCN-interior-node as a result of an indication from the threshold-metering function [RFC5670]. Abbreviated to ThM codepoint.

Excess-traffic-marked codepoint: a codepoint that indicates packets that have been marked at a PCN-interior-node as a result of an indication from the excess-traffic-metering function [RFC5670]. Abbreviated to ETM codepoint.

Not-marked codepoint: a codepoint that indicates PCN-packets that are not PCN-marked. Abbreviated to NM codepoint.

Not-PCN codepoint: a codepoint that indicates packets that are not PCN-packets.

2.2. List of Abbreviations

The following abbreviations are used in this document:

- o AF = Assured Forwarding [RFC2597]
- o CE = Congestion Experienced [RFC3168]
- o CS = Class Selector [RFC2474]
- o DSCP = Diffserv codepoint
- o e2e = end-to-end
- o ECN = Explicit Congestion Notification [RFC3168]
- o ECT = ECN Capable Transport [RFC3168]
- o EF = Expedited Forwarding [RFC3246]
- o ETM = excess-traffic-marked
- o EXP = Experimental
- o NM = not-marked
- o PCN = Pre-Congestion Notification
- o PHB = Per-hop behaviour [RFC2474]
- o ThM = threshold-marked

3. Definition of 3-in-1 PCN Encoding

The 3-in-1 PCN encoding scheme supports networks that need three PCN-marking states to be encoded within the IP header, as well as those that need only two. The full encoding is shown in Figure 1.

DSCP	Codepoint in ECN field of IP header <RFC3168 codepoint name>			
	00 <Not-ECT>	10 <ECT(0)>	01 <ECT(1)>	11 <CE>
DSCP n	not-PCN	NM	ThM	ETM

Figure 1: 3-in-1 PCN Encoding

A PCN-node will be configured to recognise certain DSCPs as PCN-compatible. Appendix A discusses the choice of suitable DSCPs. In Figure 1, 'DSCP n' indicates such a PCN-compatible DSCP. In the PCN-domain, any packet carrying a PCN-compatible DSCP and with the ECN-field anything other than 00 (not-PCN) is a PCN-packet as defined in [RFC5559].

PCN-nodes MUST interpret the ECN field of a PCN-packet using the 3-in-1 PCN encoding, rather than [RFC3168]. This does not change the behaviour for any packet with a DSCP that is not PCN-compatible, or for any node outside a PCN-domain. In all such cases, the 3-in-1 encoding is not applicable, and so by default the node will interpret the ECN field using [RFC3168].

When using the 3-in-1 encoding, the codepoints of the ECN field have the following meanings:

not-PCN: indicates a non-PCN-packet, i.e., a packet that uses a PCN-compatible DSCP but is not subject to PCN metering and marking.

NM: not-marked. Indicates a PCN-packet that has not yet been marked by any PCN marker.

ThM: threshold-marked. Indicates a PCN-packet that has been marked by a threshold-marker [RFC5670].

ETM: excess-traffic-marked. Indicates a PCN-packet that has been marked by an excess-traffic-marker [RFC5670].

4. Requirements for and Applicability of 3-in-1 PCN Encoding

4.1. PCN Requirements

In accordance with the PCN architecture [RFC5559], PCN-ingress-nodes control packets entering a PCN-domain. Packets belonging to PCN-controlled flows are subject to PCN-metering and PCN-marking, and PCN-ingress-nodes mark them as not-marked (PCN-colouring). All nodes in the PCN-domain perform PCN-metering and PCN-mark PCN-packets if needed. There are two different metering and marking behaviours: threshold-marking and excess-traffic-marking [RFC5670]. Some edge behaviours require only a Single Marking behaviour [RFC6662], others require both [RFC6661]. In the latter case, three PCN marking states are needed: not-marked (NM) to indicate not-marked packets, threshold-marked (ThM) to indicate packets marked by the threshold-marker, and excess-traffic-marked (ETM) to indicate packets marked by the excess-traffic-marker [RFC5670]. Threshold-marking and excess-traffic-marking are configured to start marking packets at different load conditions, so one marking behaviour indicates more severe pre-congestion than the other. Therefore, a fourth PCN marking state indicating that a packet is marked by both markers is not needed. However, a fourth codepoint is required to indicate packets that use a PCN-compatible DSCP but do not use PCN-marking (the not-PCN codepoint).

In all current PCN edge behaviours that use two marking behaviours [RFC5559] [RFC6661], excess-traffic-marking is configured with a larger reference rate than threshold-marking. We take this as a rule and define excess-traffic-marked as a more severe PCN-mark than threshold-marked.

4.2. Requirements Imposed by Tunnelling

[RFC6040] defines rules for the encapsulation and decapsulation of ECN markings within IP-in-IP tunnels. The publication of RFC 6040 removed the tunnelling constraints that existed when the encoding of [RFC5696] was written (see Section 3.3.2 of [RFC6627]).

Nonetheless, there is still a problem if there are any legacy (pre-RFC6040) decapsulating tunnel endpoints within a PCN-domain. If a PCN-node Threshold-marks the outer header of a tunnelled packet that has a not-marked codepoint on the inner header, a legacy decapsulator will forward the packet as not-marked, losing the Threshold-marking. The rules on applicability in Section 4.3 below are designed to avoid this problem.

Even if an operator accidentally breaks these applicability rules, the order of severity of the 3-in-1 codepoints was chosen to protect other PCN or non-PCN traffic. Although legacy pre-RFC6040 tunnels did not propagate '01', all tunnels pre-RFC6040 and post-RFC6040 have always propagated '11' correctly. Therefore, '11' was chosen to signal the most severe pre-congestion (ETM), so it would act as a reliable fail-safe even if an overlooked legacy tunnel was suppressing '01' (ThM) signals.

4.3. Applicable Environments for the 3-in-1 PCN Encoding

The 3-in-1 encoding is applicable in situations where two marking behaviours are being used in the PCN-domain. The 3-in-1 encoding can also be used with only one marking behaviour, in which case one of the codepoints MUST NOT be used anywhere in the PCN-domain (see Section 5.2.3).

With one exception (see next paragraph), any tunnel endpoints (IP-in-IP and IPsec) within the PCN-domain MUST comply with the ECN encapsulation and decapsulation rules set out in [RFC6040] (see Section 4.2).

Operators may not be able to upgrade every pre-RFC6040 tunnel endpoint within a PCN-domain. In such circumstances, a limited version of the 3-in-1 encoding can still be used but only under the following stringent condition. If any pre-RFC6040 tunnel decapsulator exists within a PCN-domain, then every PCN-node in the PCN-domain MUST be configured so that it never sets the ThM codepoint. PCN-interior-nodes in this case MUST solely use the Excess-traffic-marking function, as defined in Section 5.2.3.1. In all other situations where legacy tunnel decapsulators might be present within the PCN-domain, the 3-in-1 encoding is not applicable.

5. Behaviour of a PCN-node to Comply with the 3-in-1 PCN Encoding

Any tunnel endpoint implementation on a PCN-node MUST comply with [RFC6040]. Since PCN is a new capability, this is considered a reasonable requirement.

5.1. PCN-Ingress-Node Behaviour

If packets arrive from another Diffserv domain, any re-mapping of Diffserv codepoints MUST happen before PCN-ingress processing.

At each logical ingress link into a PCN-domain, each PCN-ingress-node will apply the four functions described in Section 4.2 of [RFC5559] to arriving packets. These functions are applied in the following order: PCN-classify, PCN-police, PCN-colour, PCN-rate-meter. This

section describes these four steps, but only the aspects relevant to packet encoding:

1. PCN-classification: The PCN-ingress-node determines whether each packet matches the filter spec of an admitted flow. Packets that match are defined as PCN-packets.
- 1b. Extra step if ECN and PCN coexist: If a packet classified as a PCN-packet arrives with the ECN field already set to a value other than Not-ECT (i.e., it is from an ECN-capable transport), then to comply with BCP 124 [RFC4774] it MUST pass through one of the following preparatory steps before the PCN-policing and PCN-colouring steps. The choice between these four actions depends on local policy:
 - * Encapsulate ECN-capable PCN-packets across the PCN-domain:
 - + either within another IP header using an RFC6040 tunnel;
 - + or within a lower-layer protocol capable of being PCN-marked, such as MPLS (see Appendix C).

Encapsulation using either of these methods is the RECOMMENDED policy for ECN-capable PCN-packets, and implementations SHOULD use IP-in-IP tunnelling as the default.

If encapsulation is used, it MUST precede PCN-policing and PCN-colouring so that the encapsulator and decapsulator are logically outside the PCN-domain (see Appendix B and specifically Figure 2).

If MPLS encapsulation is used, note that penultimate hop popping [RFC3031] is incompatible with PCN, unless the penultimate hop applies the PCN-egress-node behaviour before it pops the PCN-capable MPLS label.

- * If some form of encapsulation is not possible, the PCN-ingress-node can allow through ECN-capable packets without encapsulation, but it MUST drop CE-marked packets at this stage. Failure to drop CE-marked packets would risk congestion collapse, because without encapsulation there is no mechanism to propagate the CE markings across the PCN-domain (see Appendix B).

This policy is NOT RECOMMENDED because there is no tunnel to protect the e2e ECN capability, which is otherwise disabled when the PCN-egress-node zeroes the ECN field.

- * Drop the packet.

This policy is also NOT RECOMMENDED, because it precludes the possibility that e2e ECN can coexist with PCN as a means of controlling congestion.

- * Any other action that complies with [RFC4774] (see Appendix B for an example).

Appendix B provides more information about the coexistence of PCN and ECN.

2. PCN-policing: The PCN-policing function only allows appropriate packets into the PCN behaviour aggregate. Per-flow policing actions may be required to block rejected flows and to rate-police accepted flows, but these are specified in the relevant edge-behaviour document, e.g., [RFC6662] or [RFC6661].

Here, we only specify packet-level PCN-policing, which prevents packets that are not PCN-packets from being forwarded into the PCN-domain if PCN-interior-nodes would otherwise mistake them for PCN-packets. A non-PCN-packet will be confused with a PCN-packet if on arrival it meets all three of the following conditions:

- a) it is not classified as a PCN-packet;
- b) it already carries a PCN-compatible DSCP; and
- c) its ECN field carries a codepoint other than Not-ECT.

The PCN-ingress-node MUST police packets that meet all three conditions (a-c) by subjecting them to one of the following treatments:

- * re-mark the DSCP to a DSCP that is not PCN-compatible;
- * tunnel the packet to the PCN-egress-node with a DSCP in the outer header that is not PCN-compatible; or
- * drop the packet (NOT RECOMMENDED -- see below).

The choice between these actions depends on local policy. In the absence of any operator-specific configuration for this case, an implementation SHOULD re-mark the DSCP to zero (000000) by default.

Whichever policing action is chosen, the PCN-ingress-node SHOULD log the event and MAY also raise an alarm. Alarms SHOULD be rate-limited so that the anomalous packets will not amplify into a flood of alarm messages.

Rationale: Traffic that meets all three of the above conditions (a-c) is not PCN-traffic; therefore, ideally a PCN-ingress ought not to interfere with it, but it has to do something to avoid ambiguous packet markings. Clearing the ECN field is not an appropriate policing action, because a network node ought not to interfere with an e2e signal. Even if such packets seem like an attack, drop would be overkill, because such an attack can be neutralised by just re-marking the DSCP. And DSCP re-marking in the network is legitimate, because the DSCP is not considered an e2e signal.

3. PCN-colouring: If a packet has been classified as a PCN-packet, once it has been policed, the PCN-ingress-node:
 - * MUST set a PCN-compatible Diffserv codepoint on all PCN-packets. To conserve DSCPs, DSCPs SHOULD be chosen that are already defined for use with admission-controlled traffic. Appendix A gives guidance to implementors on suitable DSCPs.
 - * MUST set the PCN codepoint of all PCN-packets to not-marked (NM).
4. PCN rate-metering: This fourth step may be necessary depending on the edge behaviour in force. It is listed for completeness, but it is not relevant to this encoding document.

5.2. PCN-Interior-Node Behaviour

5.2.1. Behaviour Common to All PCN-Interior-Nodes

Interior nodes MUST NOT change not-PCN to any other codepoint.

Interior nodes MUST NOT change NM to not-PCN.

Interior nodes MUST NOT change ThM to NM or not-PCN.

Interior nodes MUST NOT change ETM to any other codepoint.

5.2.2. Behaviour of PCN-Interior-Nodes Using Two PCN-Markings

If the threshold-meter function indicates a need to mark a packet, the PCN-interior-node MUST change NM to ThM.

If the excess-traffic-meter function indicates a need to mark a packet:

- o the PCN-interior-node MUST change NM to ETM;
- o the PCN-interior-node MUST change ThM to ETM.

If both the threshold meter and the excess-traffic meter indicate the need to mark a packet, the Excess-traffic-marking rules MUST take precedence.

5.2.3. Behaviour of PCN-Interior-Nodes Using One PCN-Marking

Some PCN edge behaviours require only one PCN-marking within the PCN-domain. The Single Marking edge behaviour [RFC6662] requires PCN-interior-nodes to mark packets using the excess-traffic-meter function [RFC5670]. It is possible that future schemes may require only the threshold-meter function. Appendix D explains the rationale for the behaviours defined in this section.

5.2.3.1. Marking Using Only the Excess-Traffic-Meter Function

The threshold-traffic-meter function SHOULD be disabled and MUST NOT trigger any packet marking.

The PCN-interior-node SHOULD raise a management alarm if it receives a ThM packet, but the frequency of such alarms SHOULD be limited.

If the excess-traffic-meter function indicates a need to mark the packet:

- o the PCN-interior-node MUST change NM to ETM;
- o the PCN-interior-node MUST change ThM to ETM. It SHOULD also raise an alarm as above.

5.2.3.2. Marking Using Only the Threshold-Meter Function

The excess-traffic-meter function SHOULD be disabled and MUST NOT trigger any packet marking.

The PCN-interior-node SHOULD raise a management alarm if it receives an ETM packet, but the frequency of such alarms SHOULD be limited.

If the threshold-meter function indicates a need to mark the packet:

- o the PCN-interior-node MUST change NM to ThM;

- o the PCN-interior-node MUST NOT change ETM to any other codepoint. It SHOULD raise an alarm as above if it encounters an ETM packet.

5.3. PCN-Egress-Node Behaviour

A PCN-egress-node SHOULD set the not-PCN ('00') codepoint on all packets it forwards out of the PCN-domain.

The only exception to this is if the PCN-egress-node is certain that revealing other codepoints outside the PCN-domain won't contravene the guidance given in [RFC4774]. For instance, if the PCN-ingress-node has explicitly informed the PCN-egress-node that this flow is ECN-capable, then it might be safe to expose other ECN codepoints. Appendix B gives details of how such schemes might work, but such schemes are currently only tentative ideas.

If the PCN-domain is configured to use only Excess-traffic-marking, the PCN-egress-node MUST treat ThM as ETM; if only threshold-marking is used, it SHOULD treat ETM as ThM. However, it SHOULD raise a management alarm in either case since this means there is some misconfiguration in the PCN-domain.

6. Backward Compatibility

6.1. Backward Compatibility with ECN

BCP 124 [RFC4774] gives guidelines for specifying alternative semantics for the ECN field. It sets out a number of factors to be taken into consideration. It also suggests various techniques to allow the coexistence of default ECN and alternative ECN semantics. The encoding specified in this document uses one of those techniques; it defines PCN-compatible Diffserv codepoints as no longer supporting the default ECN semantics within a PCN-domain. As such, this document is compatible with BCP 124.

There is not enough space in one IP header for the 3-in-1 encoding to support both ECN marking end-to-end and PCN-marking within a PCN-domain. The non-normative Appendix B discusses possible ways to do this, e.g., by carrying e2e ECN across a PCN-domain within the inner header of an IP-in-IP tunnel. The normative text in Section 5.1 requires one of these methods to be configured at the PCN-ingress-node and recommends that implementations offer tunnelling as the default.

In any PCN deployment, traffic can only enter the PCN-domain through PCN-ingress-nodes and leave through PCN-egress-nodes. PCN-ingress-nodes ensure that any packets entering the PCN-domain have the ECN field in their outermost IP header set to the appropriate codepoint.

PCN-egress-nodes then guarantee that the ECN field of any packet leaving the PCN-domain has appropriate ECN semantics. This prevents unintended leakage of ECN marks into or out of the PCN-domain, and thus reduces backward-compatibility issues.

6.2. Backward Compatibility with the Encoding in RFC 5696

Section 5.1 of the PCN architecture [RFC5559] gives general guidance on fault detection and diagnosis, including management analysis of PCN markings arriving at PCN-egress-nodes to detect early signs of potential faults. Because the PCN encoding has gone through an obsoleted earlier stage [RFC5696], misconfiguration mistakes may be more likely. Therefore, extra monitoring, such as in the following example, may be necessary to detect and diagnose potential problems:

Informational example: In a controlled-load edge-behaviour scenario it could be worth the PCN-egress-node detecting the onset of excess-traffic marking without any prior threshold-marking. This might indicate that an interior node has been wrongly configured to mark only ETM (which would have been correct for the single-marking edge behaviour).

A PCN-node implemented to use the obsoleted encoding in RFC 5696 could conceivably have been configured so that the Threshold-meter function marked what is now defined as the ETM codepoint in the 3-in-1 encoding. However, there is no known deployment of this rather unlikely variant of RFC 5696 and no reason to believe that such an implementation would ever have been built. Therefore, it seems safe to ignore this issue.

7. Security Considerations

PCN-marking only carries a meaning within the confines of a PCN-domain. This encoding document is intended to stand independently of the architecture used to determine how specific packets are authorised to be PCN-marked, which will be described in separate documents on PCN-boundary-node behaviour.

This document assumes the PCN-domain to be entirely under the control of a single operator, or a set of operators who trust each other. However, future extensions to PCN might include inter-domain versions where trust cannot be assumed between domains. If such schemes are proposed, they must ensure that they can operate securely despite the lack of trust. However, such considerations are beyond the scope of this document.

One potential security concern is the injection of spurious PCN-marks into the PCN-domain. However, these can only enter the domain if a PCN-ingress-node is misconfigured. The precise impact of any such misconfiguration will depend on which of the proposed PCN-boundary-node behaviours is used; however, in general, spurious marks will lead to admitting fewer flows into the domain or potentially terminating too many flows. In either case, good management should be able to quickly spot the problem since the overall utilisation of the domain will rapidly fall.

8. Conclusions

The 3-in-1 PCN encoding uses a PCN-compatible DSCP and the ECN field to encode PCN-marks. One codepoint allows non-PCN traffic to be carried with the same PCN-compatible DSCP and three other codepoints support three PCN marking states with different levels of severity. In general, the use of this PCN encoding scheme presupposes that any tunnel endpoints within the PCN-domain comply with [RFC6040].

9. Acknowledgements

Many thanks to Philip Eardley for providing extensive feedback, criticism and advice. Thanks also to Teco Boot, Kwok Ho Chan, Ruediger Geib, Georgios Karagiannis, James Polk, Tom Taylor, Adrian Farrel, and everyone else who has commented on the document.

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.
- [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.
- [RFC5559] Eardley, P., "Pre-Congestion Notification (PCN) Architecture", RFC 5559, June 2009.
- [RFC5670] Eardley, P., "Metering and Marking Behaviour of PCN-Nodes", RFC 5670, November 2009.

[RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, November 2010.

10.2. Informative References

- [RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", RFC 2597, June 1999.
- [RFC3031] Rosen, E., Viswanathan, A., and R. Callon, "Multiprotocol Label Switching Architecture", RFC 3031, January 2001.
- [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.
- [RFC3540] Spring, N., Wetherall, D., and D. Ely, "Robust Explicit Congestion Notification (ECN) Signaling with Nonces", RFC 3540, June 2003.
- [RFC4594] Babiarez, J., Chan, K., and F. Baker, "Configuration Guidelines for DiffServ Service Classes", RFC 4594, August 2006.
- [RFC4774] Floyd, S., "Specifying Alternate Semantics for the Explicit Congestion Notification (ECN) Field", BCP 124, RFC 4774, November 2006.
- [RFC5127] Chan, K., Babiarez, J., and F. Baker, "Aggregation of Diffserv Service Classes", RFC 5127, February 2008.
- [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.
- [RFC5696] Moncaster, T., Briscoe, B., and M. Menth, "Baseline Encoding and Transport of Pre-Congestion Information", RFC 5696, November 2009.
- [RFC5865] Baker, F., Polk, J., and M. Dolly, "A Differentiated Services Code Point (DSCP) for Capacity-Admitted Traffic", RFC 5865, May 2010.

- [RFC6627] Karagiannis, G., Chan, K., Moncaster, T., Menth, M., Eardley, P., and B. Briscoe, "Overview of Pre-Congestion Notification Encoding", RFC 6627, July 2012.
- [RFC6661] Charny, A., Huang, F., Karagiannis, G., Menth, M., and T. Taylor, Ed., "Pre-Congestion Notification (PCN) Boundary-Node Behaviour for the Controlled Load (CL) Mode of Operation", RFC 6661, July 2012.
- [RFC6662] Charny, A., Zhang, J., Karagiannis, G., Menth, M., and T. Taylor, Ed., "Pre-Congestion Notification (PCN) Boundary-Node Behaviour for the Single Marking (SM) Mode of Operation", RFC 6662, July 2012.

Appendix A. Choice of Suitable DSCPs

This appendix is informative not normative.

A single DSCP has not been defined for use with PCN for several reasons. Firstly, the PCN mechanism is applicable to a variety of different traffic classes. Secondly, Standards Track DSCPs are in increasingly short supply. Thirdly, PCN is not a scheduling behaviour -- rather, it should be seen as being a marking behaviour similar to ECN but intended for inelastic traffic. The choice of which DSCP is most suitable for a given PCN-domain is dependent on the nature of the traffic entering that domain and the link rates of all the links making up that domain. In PCN-domains with sufficient aggregation, the appropriate DSCPs would currently be those for the Real-Time Treatment Aggregate [RFC5127]. It is suggested that admission control could be used for the following service classes (defined in [RFC4594] unless otherwise stated):

- o Telephony (EF)
- o Real-time interactive (CS4)
- o Broadcast Video (CS3)
- o Multimedia Conferencing (AF4)
- o the VOICE-ADMIT codepoint defined in [RFC5865].

CS5 is excluded from this list since PCN is not expected to be applied to signalling traffic.

PCN-marking is intended to provide a scalable admission-control mechanism for traffic with a high degree of statistical multiplexing. PCN-marking would therefore be appropriate to apply to traffic in the above classes, but only within a PCN-domain containing sufficiently aggregated traffic. In such cases, the above service classes may well all be subject to a single forwarding treatment (treatment aggregate [RFC5127]). However, this does not imply all such IP traffic would necessarily be identified by one DSCP -- each service class might keep a distinct DSCP within the highly aggregated region [RFC5127].

Guidelines for conserving DSCPs by allowing non-admission-controlled-traffic to compete with PCN-traffic are given in Appendix B.1 of [RFC5670].

Additional service classes may be defined for which admission control is appropriate, whether through some future standards action or through local use by certain operators, e.g., the Multimedia Streaming service class (AF3). This document does not preclude the use of PCN in more cases than those listed above.

Note: The above discussion is informative not normative, as operators are ultimately free to decide whether to use admission control for certain service classes and whether to use PCN as their mechanism of choice.

Appendix B. Coexistence of ECN and PCN

This appendix is informative not normative. It collects together material relevant to coexistence of ECN and PCN, including that spread throughout the body of this specification. If this results in any conflict or ambiguity, the normative text in the body of the specification takes precedence.

ECN [RFC3168] is an e2e congestion notification mechanism. As such it is possible that some traffic entering the PCN-domain may also be ECN-capable. The PCN encoding described in this document reuses the bits of the ECN field in the IP header. Consequently, this disables ECN within the PCN-domain.

For the purposes of this appendix, we define two forms of traffic that might arrive at a PCN-ingress-node. These are admission-controlled traffic (PCN-traffic) and non-admission-controlled traffic (non-PCN-traffic).

Flow signalling identifies admission-controlled traffic, by associating a filter spec with the need for admission control (e.g., through RSVP or some equivalent message, such as from a SIP server to the ingress or from a logically centralised network control system). The PCN-ingress-node re-marks admission-controlled traffic matching that filter spec to a PCN-compatible DSCP. Note that the term "flow" need not imply just one microflow, but instead could match an aggregate and/or could depend on the incoming DSCP (see Appendix A).

All other traffic can be thought of as non-admission-controlled (and therefore outside the scope of PCN). However, such traffic may still need to share the same DSCP as the admission-controlled traffic. This may be due to policy (for instance, if it is high-priority voice traffic), or may be because there is a shortage of local DSCPs.

Unless specified otherwise, for any of the cases in the list below, an IP-in-IP tunnel that complies with [RFC6040] can be used to preserve ECN markings across the PCN-domain. The tunnelling action

should be applied wholly outside the PCN-domain as illustrated in Figure 2. Then, by the rules of RFC 6040, the tunnel egress propagates the ECN field from the inner header, because the PCN-egress will have zeroed the outer ECN field.

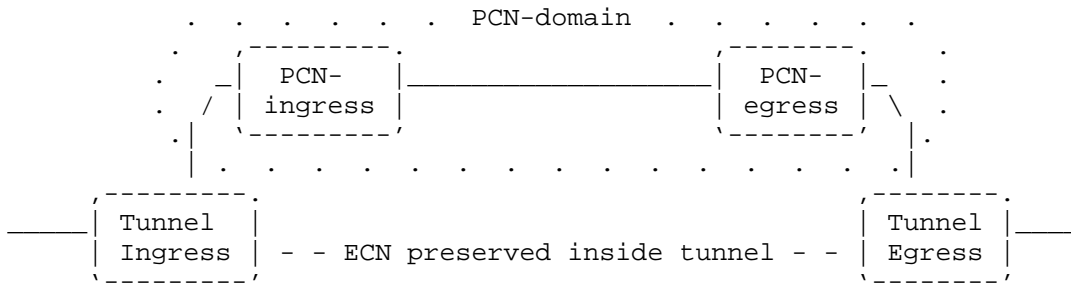


Figure 2: Separation of Tunnelling and PCN Actions

There are three cases for how e2e ECN traffic may wish to be treated while crossing a PCN-domain:

- a) Traffic that does not require PCN admission control:
 For example, traffic that does not match flow signalling being used for admission control. In this case, the e2e ECN traffic is not treated as PCN-traffic. There are two possible scenarios:
- * Arriving traffic does not carry a PCN-compatible DSCP: no action required.
 - * Arriving traffic carries a DSCP that clashes with a PCN-compatible DSCP. There are two options:
 1. The ingress maps the DSCP to a local DSCP with the same scheduling PHB as the original DSCP, and the egress re-maps it to the original PCN-compatible DSCP.
 2. The ingress tunnels the traffic, setting the DSCP in the outer header to a local DSCP with the same scheduling PHB as the original DSCP.
 3. The ingress tunnels the traffic, using the original DSCP in the outer header but setting not-PCN in the outer header; note that this turns off ECN for this traffic within the PCN-domain.

The first or second options are recommended unless the operator is short of local DSCPs.

b) Traffic that requires admission-control:

In this case, the e2e ECN traffic is treated as PCN-traffic across the PCN-domain. There are two options.

- * The PCN-ingress-node places this traffic in a tunnel with a PCN-compatible DSCP in the outer header. The PCN-egress zeroes the ECN-field before decapsulation.
- * The PCN-ingress-node drops CE-marked packets and otherwise sets the ECN-field to NM and sets the DSCP to a PCN-compatible DSCP. The PCN-egress zeroes the ECN field of all PCN packets; note that this turns off e2e ECN.

The second option is emphatically not recommended, unless perhaps as a last resort if tunnelling is not possible for some insurmountable reason.

c) Traffic that requires PCN admission control where the endpoints ask to see PCN marks:

Note that this scheme is currently only a tentative idea.

For real-time data generated by an adaptive codec, schemes have been suggested where PCN marks may be leaked out of the PCN-domain so that end hosts can drop to a lower data-rate, thus deferring the need for admission control. Currently, such schemes require further study and the following is for guidance only.

The PCN-ingress-node needs to tunnel the traffic as in Figure 2, taking care to comply with [RFC6040]. In this case, the PCN-egress does not zero the ECN field (contrary to the recommendation in Section 5.3), so that the [RFC6040] tunnel egress will preserve any PCN-marking. Note that a PCN-interior-node may change the ECN-field from '10' to '01' (NM to ThM), which would be interpreted by the e2e ECN as a change from ECT(0) into ECT(1). This would not be compatible with the (currently experimental) ECN nonce [RFC3540].

Appendix C. Example Mapping between Encoding of PCN-Marks in IP and in MPLS Shim Headers

This appendix is informative not normative.

The 6 bits of the DS field in the IP header provide for 64 codepoints. When encapsulating IP traffic in MPLS, it is useful to make the DS field information accessible in the MPLS header. However, the MPLS shim header has only a 3-bit traffic class (TC) field [RFC5462] providing for 8 codepoints. The operator has the freedom to define a site-local mapping of the 64 codepoints of the DS field onto the 8 codepoints in the TC field.

[RFC5129] describes how ECN markings in the IP header can also be mapped to codepoints in the MPLS TC field. Appendix A of [RFC5129] gives an informative description of how to support PCN in MPLS by extending the way MPLS supports ECN. [RFC5129] was written while PCN specifications were in early draft stages. The following provides a clearer example of a mapping between PCN in IP and in MPLS using the PCN terminology and concepts that have since been specified.

To support PCN in a MPLS domain, a PCN-compatible DSCP ('DSCP n') needs codepoints to be provided in the TC field for all the PCN-marks used. That means, when, for instance, only excess-traffic-marking is used for PCN purposes, the operator needs to define a site-local mapping to two codepoints in the MPLS TC field for IP headers with:

- o DSCP n and NM
- o DSCP n and ETM

If both excess-traffic-marking and threshold-marking are used, the operator needs to define a site-local mapping to codepoints in the MPLS TC field for IP headers with all three of the 3-in-1 codepoints:

- o DSCP n and NM
- o DSCP n and ThM
- o DSCP n and ETM

In either case, if the operator wishes to support the same Diffserv PHB but without PCN marking, it will also be necessary to define a site-local mapping to an MPLS TC codepoint for IP headers marked with:

- o DSCP n and not-PCN

The above sets of codepoints are required for every PCN-compatible DSCP. Clearly, given so few TC codepoints are available, it may be necessary to compromise by merging together some capabilities. Guidelines for conserving TC codepoints by allowing non-admission-controlled-traffic to compete with PCN-traffic are given in Appendix B.1 of [RFC5670].

Appendix D. Rationale for Difference between the Schemes Using One PCN-Marking

Readers may notice a difference between the two behaviours in Sections 5.2.3.1 and 5.2.3.2. With only Excess-traffic-marking enabled, an unexpected ThM packet can be re-marked to ETM. However, with only Threshold-marking, an unexpected ETM packet cannot be re-marked to ThM.

This apparent inconsistency is deliberate. The behaviour with only Threshold-marking keeps to the rule of Section 5.2.1 that ETM is more severe and must never be changed to ThM even though ETM is not a valid marking in this case. Otherwise, implementations would have to allow operators to configure an exception to this rule, which would not be safe practice and would require different code in the data plane compared to the common behaviour.

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/>

Toby Moncaster
University of Cambridge Computer Laboratory
William Gates Building, J J Thomson Avenue
Cambridge CB3 0FD
UK

EMail: toby.moncaster@cl.cam.ac.uk

Michael Menth
University of Tuebingen
Sand 13
72076 Tuebingen
Germany

Phone: +49-7071-2970505
EMail: menth@uni-tuebingen.de

