

Inter-Domain Routing Working Group
Internet-Draft
Intended status: Standards Track
Expires: May 23, 2020

Th. Knoll
November 20, 2019

BGP Class of Service Interconnection
draft-knoll-idr-cos-interconnect-23

Abstract

This document focuses on Class of Service Interconnection at inter-domain interconnection points. It specifies two new transitive attributes, which enable adjacent peers to signal Class of Service Capabilities and certain Class of Service admission control Parameters. The new "CoS Capability" is deliberately kept simple and denotes the general EF, AF Group BE and LE forwarding support across the advertising AS. The second "CoS Parameter Attribute" is of variable length and contains a more detailed description of available forwarding behaviours using the PHB id Code encoding. Each PHB id Code is associated with rate and size based traffic parameters, which will be applied in the ingress AS Border Router for admission control purposes to a given forwarding behaviour.

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].

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 May 23, 2020.

Copyright Notice

Copyright (c) 2019 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	2
2. Definition and Usage of the CoS Capability	3
2.1. Extended Community Type	3
2.2. Structure of the CoS Capability Attribute	4
2.3. Usage of the CoS Capability Attribute	7
3. Definition and Usage of the CoS Parameter Attribute	7
3.1. Definition of the CoS Parameter Attribute	7
3.2. Usage of the CoS Parameter Attribute	8
4. Confidentiality Considerations	9
5. IANA Considerations	9
6. Security Considerations	10
7. References	10
7.1. Normative References	10
7.2. Informative References	11
Author's Address	11

1. Introduction

AS interconnection is currently based on best effort interconnection only. BGP-4 [RFC4271] is the de-facto interconnection protocol used to exchange reachability information. There is no standardized set of supported traffic classes, no standardized packet marking and no standardized forwarding behaviour, which cross-domain traffic could rely on. QoS policy decisions are taken by AS providers independently and in an uncoordinated fashion. However, many AS providers make use of the Differentiated Services Architecture [RFC2475] as AS internal QoS mechanism. Within this architecture, there are 64 codepoints and an unlimited number of Per Hop Behaviours (PHBs) available. Some PHBs have been defined in separate RFCs, which will be focused on in this document.

A Basic Set of supported Classes, called "Basic CoS" is defined inhere, which consists of the primitive "Best Effort (BE)" PHB, the "Expedited Forwarding (EF)" PHB [RFC3246], the "Assured Forwarding (AF)" PHB Group [RFC2597] and the "Lower Effort" Per-Domain Behavior (PDB) [RFC3662]. AS providers, which can support this Basic CoS are asked to signal this capability to their interconnection partners by means of the new CoS Capability Extended Community defined in Section 2 of this draft.

4 AF PHB classes have been defined so far, which will be grouped into the generally signalled "AF Group". That is, as long as the AS provider can support at least one out of the 4 AF classes in his externally supported CoS Set, this AS is regarded as AF capable.

A second transitive attribute is defined in Section 3, which is used for parameter signalling about the applied access control within the ingress AS border router. The reason for this traffic limitation is the fact, that certain high quality forwarding behaviours can only be achieved, if the percentage of high priority traffic within the traffic mix lies below a certain threshold. This attribute informs the interconnection partner about the applied limitation, which can in turn be used to perform traffic shaping at the neighbouring AS' egress. The attribute allows this limitation signalling either associated to the NLRI within the same UPDATE message or with "global" scope to describe the generally applied ingress limitation.

Both attributes are likely to be used together, if ingress class limitation is used for the respective AS.

More detailed signalling of forwarding behaviour distinction and associated cross-layer marking can be achieved using the QoS Marking Attribute approach [I-D.knoll-idr-qos-attribute].

2. Definition and Usage of the CoS Capability

2.1. Extended Community Type

The new CoS Capability is encoded as a BGP Extended Community [RFC4360]. Extended Community Attributes are transitive optional BGP attributes with Type Code 16. An adoption to the simple BGP Community Attribute encoding [RFC1997] is not defined in this document. The actual encoding within the BGP Extended Community Attribute is as follows.

The CoS Capability is transitive and of regular type which results in a 1 octet Type field followed by 7 octets for the CoS Capability structure. The Type is IANA-assignable (FCFS procedure) and marks

the community as transitive across ASes. The type number has been assigned by IANA to 0xYY (0x00-0x3f).

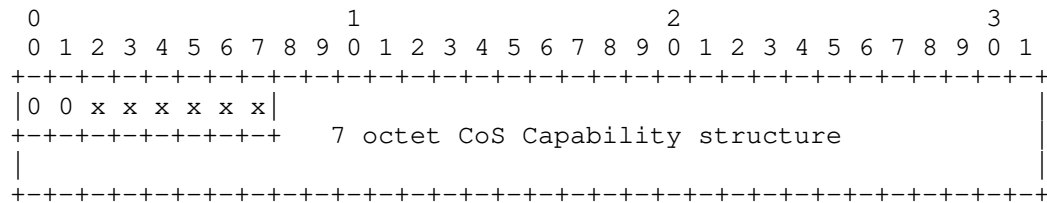


Figure 1

2.2. Structure of the CoS Capability Attribute

The CoS Capability structure is deliberately kept very simple and is defined as follows.

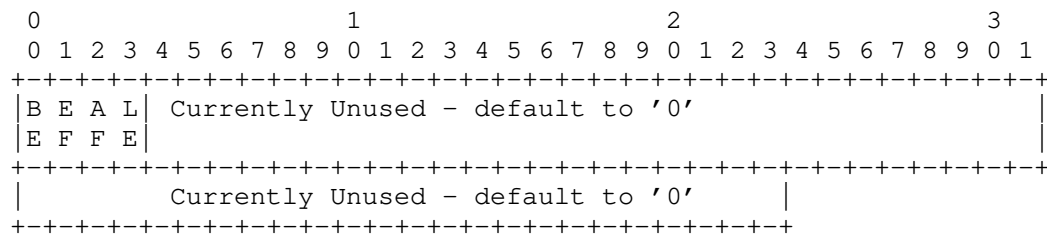


Figure 2

The Currently Unused bits default to '0' and MUST be ignored on reception.

Leading "BE, EF, AF and LE" encoding.

This encoding signals the BE, EF, AF Group and LE support of the respective AS.

Bit	Encoding
BE	Default to '1' to signal general "Best Effort" PHB support
EF	'1' ... "Expedited Forwarding" PHB support [RFC3246]
AF	'1' ... "Assured Forwarding" PHB group support [RFC2597]
LE	'1' ... "Lower Effort" PDB support [RFC3662]

Table 1: CoS support encoding

The implied Per-Hop-Behaviour Identification Codes follow the definition as standardized in [RFC3140]. The AF Group needs to consist of at least one of the currently available AF1x, AF2x, AF3x and AF4x.

BE:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
EF:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0
AF1x:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	1	0	1	0	0	0	0	0	0	0	0	0	1	0
AF2x:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	1	0	0	1	0	0	0	0	0	0	0	0	0	1	0
AF3x:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	1	1	0	1	0	0	0	0	0	0	0	0	0	1	0
AF4x:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0
LE:															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 3

2.3. Usage of the CoS Capability Attribute

The CoS Capability is used as primitive means to signal the general availability of the set of "Basic CoS" PHBs in the advertising AS. This Extended Community is included within the attribute section of an BGP UPDATE message and is therefore associated to the NLRI information of the same message. Whether the Basic CoS is available and is therefore advertised can easily being judged on for all prefixes, which originate from the advertising AS.

All other reachability information MUST be signalled together with this CoS Capability if they were received together with such an Extended Community by neighbouring peers.

NLRI MUST NOT be marked as supporting "Basic CoS" by means of the CoS Capability, if it were not received together with such an attribute.

3. Definition and Usage of the CoS Parameter Attribute

3.1. Definition of the CoS Parameter Attribute

The CoS Parameter Attribute is an optional transitive BGP attribute.

The attribute contains one or more of the following:

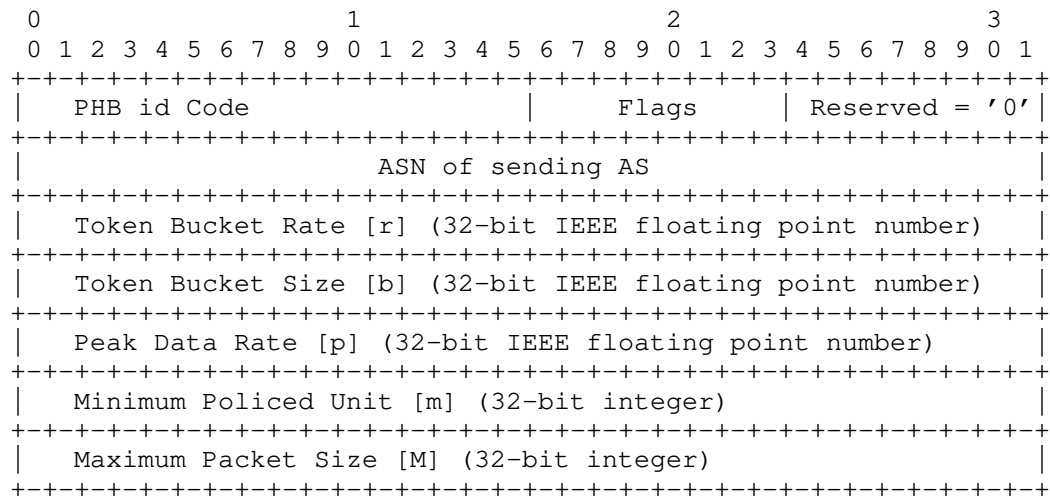


Figure 4

PHB ID:

This field specifies the targeted Per Hop Behaviour limitations and follows the defined encoding of [RFC3140] as listed in Figure 3.

Flags:

```

  0  1  2  3  4  5  6  7
+---+---+---+---+---+---+---+
| G | DR | 0 | 0 | 0 | 0 | 0 | 0 |
+---+---+---+---+---+---+---+

```

Only two flags are defined. The remaining bits default to '0' and MUST be ignored on reception.

The 'G' flag signals, whether the limitations have global scope on all incoming traffic ('1') or are associated to traffic that is destined to destinations within the NLRI of the UPDATE message ('0'). NLRI specific limitation will supersede globally signalled ones for traffic destined to those NLRI destinations.

The 'DR' flag signals the applied handling of non-confirming traffic. DR='0' signals strict dropping of excess traffic. DR='1' signals the performed remarking of excess traffic packets to Best Effort traffic marking.

ASN of sending AS:

Depending on the 2-octet or 4-octet AS peering type, the sending AS of the attribute MUST encode its AS number as right-aligned 32bit number.

Peak Data Rate, Token Bucket Rate, Token Bucket Size, Minimum Policed Unit and Maximum Packet Size:

The rates and sizes are given in 4 octet IEEE floating point format [IEEE] or 4 octet integer format, respectively. They are parameters to a token bucket ingress filter, which is applied to the packets belonging to the stated PHB id. The parameters follow the definition given in [RFC2210] and [RFC2215].

3.2. Usage of the CoS Parameter Attribute

The signalled parameters are used for PHB id Code based ingress limitation. Depending on which PHB id Codes a BGP peer signals in this attribute to its neighbour, it is said, that the respective PHB id Code is supported and will experience the defined limitations.

Those limitations can be applied to all incoming traffic of a specific PHB id Code (marked as 'G') or only for incoming traffic, that is destined for the NLRI of the given UPDATE message.

The resulting treatment for non-confirming traffic is signalled through the 'DR' flag.

To withdraw a previously signalled limitation, a CoS Parameter Attribute for the respective PHB id Code MUST be sent with a rate value [r] of zero. Using the 'G' flag, this can be withdrawn globally for all traffic of the given PHB id Code or withdrawn only for traffic destined to the prefixes given in the NLRI of the UPDATE. Previously signalled non-global (i.e. NLRI specific) limitations are also waived, if the same prefix is advertised without a CoS Parameter Attribute later on. In this case, the missing attribute is considered as the above described 'rate zero update' for those prefixes. Waived prefix specific limitations do not supersede global limitations for the respective PHB id Code. In turn, a withdrawal of a global limitation does also withdraw any possibly existing prefix specific ones for the respective PHB id Code.

All limitations have AS local scope for the advertising AS and the neighbouring AS might or might not adopt its sending behaviour to those advertised limitations.

Despite the transitive nature of the new attribute, its usage for ingress limitation is confined to neighbouring ASes. Processing of the conveyed parameters is only valid for peers, who are peering with the AS specified in the ASN field of the attribute.

The attribute SHOULD NOT be transitively relayed to non-adjacent interconnection partners.

4. Confidentiality Considerations

The disclosure of confidential AS intrinsic information by means of the signalled Basic CoS support is of low key security concern. The disclosure of information through CoS Parameter signalling is more detailed. However, all included parameters are exchanged with direct interconnection partners and are the free choice of each AS provider.

5. IANA Considerations

This document defines a new BGP Extended Community, which needs to be assigned a number by IANA within the Extended Community list. The new CoS Capability is a BGP Extended Community of regular type. It is IANA-assignable (FCFS procedure) and is transitive across ASes. A

number assignment application within the numbering range of 0x00-0x3f is made to IANA.

Note to RFC Editor: this section may be removed on publication as an RFC.

This document defines a new BGP attribute. This attribute is optional and transitive.

6. Security Considerations

This extension to BGP does not change the underlying security issues inherent in the existing BGP version.

The signalled attributes are transitive with limited relay operation in the CoS Parameter Attribute case. AS peers, which use egress traffic shaper on the signalled limitations SHOULD exhaust all available BGP security features to make sure, that the signalled limitation is actually sent by the adjacent peer.

7. References

7.1. Normative References

- [IEEE] IEEE, "IEEE Standard for Binary Floating-Point Arithmetic", ISBN 1-5593-7653-8, 1985.
- [RFC1997] Chandra, R., Traina, P., and T. Li, "BGP Communities Attribute", RFC 1997, DOI 10.17487/RFC1997, August 1996, <<https://www.rfc-editor.org/info/rfc1997>>.
- [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>>.
- [RFC2210] Wroclawski, J., "The Use of RSVP with IETF Integrated Services", RFC 2210, DOI 10.17487/RFC2210, September 1997, <<https://www.rfc-editor.org/info/rfc2210>>.
- [RFC2215] Shenker, S. and J. Wroclawski, "General Characterization Parameters for Integrated Service Network Elements", RFC 2215, DOI 10.17487/RFC2215, September 1997, <<https://www.rfc-editor.org/info/rfc2215>>.

- [RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", RFC 2597, DOI 10.17487/RFC2597, June 1999, <<https://www.rfc-editor.org/info/rfc2597>>.
- [RFC3140] Black, D., Brim, S., Carpenter, B., and F. Le Faucheur, "Per Hop Behavior Identification Codes", RFC 3140, DOI 10.17487/RFC3140, June 2001, <<https://www.rfc-editor.org/info/rfc3140>>.
- [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, DOI 10.17487/RFC3246, March 2002, <<https://www.rfc-editor.org/info/rfc3246>>.
- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, DOI 10.17487/RFC4271, January 2006, <<https://www.rfc-editor.org/info/rfc4271>>.
- [RFC4360] Sangli, S., Tappan, D., and Y. Rekhter, "BGP Extended Communities Attribute", RFC 4360, DOI 10.17487/RFC4360, February 2006, <<https://www.rfc-editor.org/info/rfc4360>>.

7.2. Informative References

- [I-D.knoll-idr-qos-attribute]
Knoll, T., "BGP Extended Community for QoS Marking", draft-knoll-idr-qos-attribute-24 (work in progress), July 2019.
- [RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and W. Weiss, "An Architecture for Differentiated Services", RFC 2475, DOI 10.17487/RFC2475, December 1998, <<https://www.rfc-editor.org/info/rfc2475>>.
- [RFC3662] Bless, R., Nichols, K., and K. Wehrle, "A Lower Effort Per-Domain Behavior (PDB) for Differentiated Services", RFC 3662, DOI 10.17487/RFC3662, December 2003, <<https://www.rfc-editor.org/info/rfc3662>>.

Author's Address

Thomas Martin Knoll

Email: thomas.m.knoll@gmail.com

Inter-Domain Routing Working Group
Internet-Draft
Intended status: Standards Track
Expires: January 26, 2020

Th. Knoll
July 25, 2019

BGP Extended Community for QoS Marking
draft-knoll-idr-qos-attribute-24

Abstract

This document specifies a simple signalling mechanism for inter-domain QoS marking using several instances of a new BGP Extended Community. Class based packet marking and forwarding is currently performed independently within ASes. The new QoS marking community makes the targeted Per Hop Behaviour within the IP prefix advertising AS and the currently applied marking at the interconnection point known to all access and transit ASes. This enables individual (re-)marking and possibly forwarding treatment adaptation to the original QoS class setup of the respective originating AS. The extended community provides the means to signal QoS markings on different layers, which are linked together in QoS Class Sets. It provides inter-domain and cross-layer insight into the QoS class mapping of the source AS with minimal signalling traffic.

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].

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 January 26, 2020.

Copyright Notice

Copyright (c) 2019 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	2
2. Problem Statement	3
3. Related Work	5
4. Definition of the QoS Marking Community	7
4.1. Extended Community Type	7
4.2. Structure of the QoS Marking Community	7
4.3. Technology Type Enumeration	10
5. Community Usage	11
5.1. QoS Marking Example	12
5.2. AS Border Packet Forwarding	12
5.3. IP Prefix Aggregation	13
6. Confidentiality Considerations	13
7. IANA Considerations	13
8. Security Considerations	14
9. References	14
9.1. Normative References	14
9.2. Informative References	15
Appendix A. QoS Marking Example	16
Author's Address	17

1. Introduction

A new BGP Extended Community is defined in this document, which carries QoS marking information for different network layer technologies across ASes. This extended community is called "QoS Marking". This new community provides a mechanism within BGP-4 [RFC4271] for associating all advertised prefixes of the AS with its differentiated QoS Class Marking information. It allows for the consistent exchange of class encoding values between BGP peers for physical, data link and network QoS mechanisms. These labels can be used to control the distribution of this information, for the

encoding and for treatment adjustments within the AS or for other applications. One globally seen QoS Class Set per AS is required for scalability reasons. It is the AS provider's responsibility to enforce the globally signalled Set throughout the AS.

Several QoS Marking communities MAY be included in a single BGP UPDATE message. They are virtually linked together by means of an identical "QoS Set Number" field. Each QoS Marking community is encoded as 8-octet tuple, as defined in Section 4. Signalled QoS Class Sets are assumed to be valid for traffic crossing this AS. If different QoS strategies are used with an AS, its provider is responsible for consistent transport of transit traffic across this inhomogeneous domain. In all transit forwarding cases, QoS based tunnelling mechanisms are the means of choice for transparent traffic transport.

The availability of the "Best Effort" forwarding class is implied and defaults to a zero encoding on all signalled layers. It is therefore not necessary to include QoS Marking communities for the Best Effort Class as long as the default encoding is in place.

Class overload prevention can be achieved by means of the signalling described in [I-D.knoll-idr-cos-interconnect]. It is a complementary concept to limit the usage of advertised classes in a fair and square manner.

2. Problem Statement

Current inter-domain interconnection is "best effort" interconnection only. That is, traffic forwarding between ASes is without traffic class differentiation and without any forwarding guarantee. It is common for network providers to reset any IP packet class markings to zero, the best effort DSCP marking, at the AS ingress router, which eliminates any traffic differentiation. Some providers perform higher layer classification at the ingress in order to guess the forwarding requirements and to match on their AS internal QoS forwarding policy. There is no standardized set of classes, no standardized marking (class encoding) and no standardized forwarding behaviour, which cross-domain traffic could rely on. QoS policy decisions are taken by network providers independently and in an uncoordinated fashion.

This general statement does not cover the existing individual agreements, which do offer quality based interconnection with strict QoS guarantees. However, such SLA based agreements are of bilateral or multilateral nature and do not offer a means for a general "better than best effort" interconnection. This draft does not aim for making such SLA based agreements become void. On the contrary, those

agreements are expected to exist for special traffic forwarding paths with strictly guaranteed QoS.

There are many approaches, which propose proper inter-domain QoS strategies including inter-domain parameter signalling, metering, monitoring and misbehaviour detection. Such complex strategies get close to guaranteed QoS based forwarding at the expense of dynamic measurements and adjustments, of state keeping on resource usage vs. traffic load and in particular of possibly frequent inter-domain signalling.

The proposed QoS Class marking approach dissociates from the complex latter solutions and targets the general "better than best effort" interconnection in coexistence with SLA based agreements. It enables ASes to make their supported Class Sets and their encoding globally known. In other words, this support information constitutes a simple map of QoS enabled roads in transit and destination ASes.

Signalling the coarse information about the supported class set and its cross-layer encoding within the involved forwarding domains of the selected AS path removes the lack of knowledge about the over-all available traffic differentiation. AS providers are enabled to make an informed decision about supported class encodings and might adopt to them. No guarantees are offered by this "better than best effort" approach, but as much as easily possible traffic differentiation without the need for frequent inter-domain signalling and for costly ingress re-classification will be achieved.

Remarking the class encoding of customer traffic in order to match neighbouring class set encodings is reasonable at AS interconnection points. For AS internal forwarding, the encapsulation within any kind of QoS supporting tunnelling technology is highly recommended. The cross-layer signalling of QoS encoding will further ease the setup of QoS based inter-domain tunnelling.

The general confidentiality concern of disclosing AS internal policy information is addressed in Section 6. In short, network providers can signal a different class set in the QoS Marking communities to the one actually used internally. The different class sets (externally signalled vs. internally applied one) require an undisclosed strictly defined mapping at the AS borders between the two. This way, a distinction between internal and external QoS Class Sets can be achieved.

The general need for class based accounting is not addressed by this draft. MIB extensions are also required, which separate traffic variables by traffic marking. It is expected for both that existing procedures can be reused in a class based manner.

3. Related Work

A number of QoS improvement approaches have been proposed before and a selection will be briefly mentioned in this section.

Most of the approaches perform parameter signalling.

[I-D.jacquetnet-bgp-qos] defines the QOS_NLRI attribute, which is used for propagating QoS-related information associated to the NLRI (Network Layer Reachability Information) information conveyed in a BGP UPDATE message. Single so called "QoS routes" are signalled, which fulfil certain QoS requirements. Several information types are defined for the attribute, which concentrate on rate and delay type parameters.

[I-D.boucadair-qos-bgp-spec] is based on the specified QOS_NLRI attribute and introduces some modifications to it. The notion of AS-local and extended QoS classes is used, which effectively describes the local set of QoS performance parameters or their cross-domain combined result. Two groups of QoS delivery services are distinguished, where the second group concentrates on ID associated QoS parameter propagation between adjacent peers. The first group is of more interest for this draft since it concentrates on the "identifier propagation" such as the DSCP value for example. However, this signalling is specified for the information exchange between adjacent peers only and assumes the existence of extended QoS classes and offline traffic engineering functions.

Another approach is described in [I-D.liang-bgp-qos]. It associates a list of QoS metrics with each prefix by extending the existing AS_PATH attribute format. Hop-by-hop metric accumulation is performed as the AS_PATH gets extended in relaying ASes. Metrics are generically specified as a list of TLV-style attribute elements. The metrics such as bandwidth and delay are exemplary mentioned in the draft.

One contribution specialized in the signalling of Type Of Service (TOS) values which are in turn directly mapped to DSCP values in section 3.2 of the draft [I-D.zhang-idr-bgp-extcommunity-qos]. The TOS value is signalled within an Extended Community Attribute and, if it is understood correctly, will be applied to a certain route. An additional value field is used to identify, which routes belong to which signalled TOS community. Who advertises such attributes and whether they are of transitive or non-transitive type remains unspecified.

The most comprehensive analysis (although not an IETF draft) is given in [MIT_CFP]. This "Inter-provider Quality of Service" white paper examines the inter-domain QoS requirements and derives a

comprehensive approach for the introduction of at least one QoS class with guaranteed delay parameters. The implementation aspects of metering, monitoring, parameter feedback and impairment allocations are all considered in the white paper. However, QoS guarantees and parameter signalling is beyond the intention of this QoS Marking draft.

Other drafts may also be considered as related work as long as they convey QoS marking information and might be "misused" for QoS class signalling.

One example is the usage of the "Traffic Engineering Attribute" as defined in [RFC5543]. However, the attribute is non-transitive and the LSP encoding types are not generally applicable to inter-domain interconnection types. Its usage of the targeted QoS Marking signalling is not possible. The included maximum bandwidth of each of eight priority classes, could however be used in future draft extensions.

The second example is the current "Dissemination of flow specification rules" draft [RFC5575]. It defines a new BGP NLRI encoding format, which can be used to distribute traffic flow specifications. Such flow specification can also include DSCP values as type 11 in the NLRI. Furthermore, one could signal configuration actions together with the DSCP encoding, which could be used for filtering purposes or even trigger remarking and route selection with it. Such usage is not defined in the draft and can hardly be achieved because of the following reasons. The flow specification is focused on single flows, which might even be part of an aggregate. Such fine grained specification is counterproductive for the coarse grained general QoS Marking approach of this draft. The novel approach of cross-layer QoS Marking could also not be incorporated, which might be essential for future tunnelled inter-domain interconnection.

4. Definition of the QoS Marking Community

4.1. Extended Community Type

The new QoS Marking community is encoded in a BGP Extended Community Attribute [RFC4360]. It is therefore a transitive optional BGP attribute with Type Code 16. An adoption to the simple BGP Community Attribute encoding [RFC1997] is not defined in this document. The actual encoding within the BGP Extended Community Attribute is as follows.

The QoS Marking community is of regular type which results in a 1 octet Type field followed by 7 octets for the QoS marking structure. The Type is IANA-assignable and marks the community as transitive across ASes. The type number has been assigned by IANA to 0x04 [IANA_EC].

Optionally, a non-transitive Type value assignment of 0x44 is provided, which allows for the AS internal marking information exchange. The community format remains untouched for the non-transitive version.

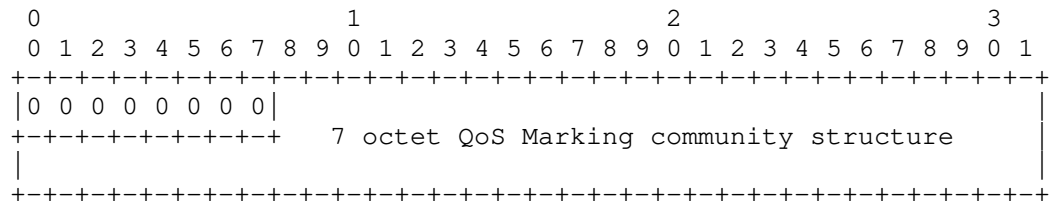


Figure 1

4.2. Structure of the QoS Marking Community

The QoS Marking community provides a flexible encoding structure for various QoS Markings on different layers. This flexibility is achieved by a Flags, a QoS Set Number and a Technology Type field within the 7 octet structure as defined below.

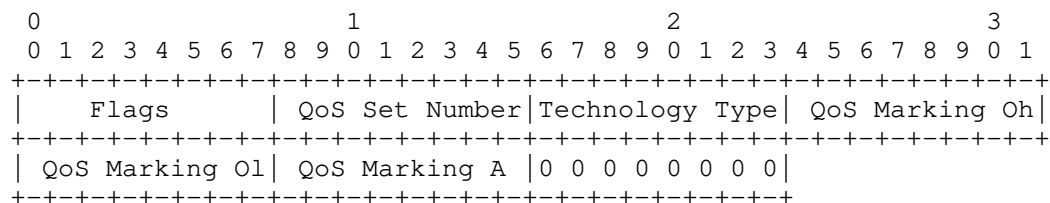


Figure 2

Flags:

```

  0  1  2  3  4  5  6  7
+--+--+--+--+--+--+--+
| 0  0 | P | R | I | A | 0  0 |
+--+--+--+--+--+--+--+

```

Figure 3

All used and unused flags default to a value of '0'. The following table shows the bit encoding of the Flags field.

Bit	Flag	Encoding
0-1	unused	Default to '0'
2	P	'1' ... Marking is preserved
3	R	'1' ... Remarking occurred
4	I	'1' ... QoS marking ignored
5	A	'1' ... QoS class aggregation occurred
6,7	unused	Default to '0'

Table 1

The 'P' flag indicates the preservation of incoming markings during the transit forwarding process. The IP prefix originating AS SHOULD set the flag to '1', which is otherwise implied by an AS_PATH length of 1 ASN. Transit ASes MUST set the flag to '1', if the advertised Marking A is accepted at the ingress and is sent out unchanged at the egress. That is, no remarking occurs - neither for marking adoption with the neighbouring downstream AS nor by resetting the markings. This flag field is set and cleared by each relaying AS according to its handling of markings - irrespective of the possible ignorance of the particular Marking A in the internal per hop forwarding behaviour.

The Flags "R, I and A" are set to '0' in the advertisement by the IP prefix originating AS. Transit ASes MUST change the flag value to '1' once the respective event occurred. If the QoS marking actively used in the transit AS internal forwarding is different from the advertised original one, the 'Remarking (R)' flag is set to '1'. This MUST be done separately for each technology type community within the community set. The same applies to the 'Ignore (I)' flag, if the respective advertised QoS marking is ignored in the transit AS internal forwarding.

The 'Aggregation (A)' flag MUST be set to '1' by the UPDATE message relaying transit AS, if the respective IP prefixes will be advertised inside an IP prefix aggregate constituted from differing Class Sets.

If the defined "R, I and A" flags are cleared - and by means of the cleared 'Partial' flag of the BGP attribute it is shown, that no "QoS Class ignorant" AS is involved in the forwarding path - a consistent class based overall traffic separated forwarding is available along this path.

QoS Set Number:

Several single QoS Marking communities can be logically grouped into a QoS Marking community Set characterized by a identical QoS Set Number. This grouping of the single QoS Marking communities into a set provides cross-layer linking between the QoS class encodings. It can also be used for the specification of behaviour sets as given in the [RFC3140]. The number of signalled QoS Marking communities as well as QoS Marking community Sets is at the operator's choice of the originating AS. The enumerated QoS set numbers have BGP UPDATE message local significance starting with set number 0x00.

Technology Type:

The technology type encoding uses the enumeration list in (Section 4.3). Future version of this draft will need an extended enumeration list administered by IANA.

QoS Marking / Enumeration O & A:

The interpretation of these fields depends on the selected layer and technology. ASes, which process the community and support the given QoS Class by means of a QoS mechanism using bit encodings for the targeted behaviour (e.g. IP DSCP, Ethernet User Priority, MPLS TC etc.) MUST use a copy of the encoding in the "QoS Marking A" community field. Unused higher order bits default to '0'. Other technologies, which use separate forwarding channels for different classes (such as L-LSPs, VPI/VCI inferred ATM classes, lambda inferred priority, etc.) SHALL use class enumerations as encoding in this community field. The enumeration count starts with zero for the best effort traffic class and rises by one with each available higher priority class.

There are two QoS Marking fields within the QoS Marking community for the "original (O)" and the "active (A)" QoS marking. Higher order

bits of those fields, which are not used for the respective behaviour encoding, default to zero.

QoS Marking O (Original QoS Marking):

This field is a 16 bit QoS Marking field, which consists of of a high ("Oh") and a low ("Ol") octet. The IP prefix originating AS copies the internally associated QoS encoding of the given Technology Type into this one octet field. The field value is right-aligned depending on the number of encoded bits. For the IP technology, the encoding of Per Hop Behaviour Codes has to follow the definitions stated in [RFC3140]. The field MUST remain unchanged in BGP UPDATE messages of relaying nodes.

QoS Marking A (Active QoS Marking):

QoS Marking A and O MUST be identically encoded by the prefix originating AS, except for the case, where IP technology Per Hop Behaviours are addressed. "QoS Marking A" will always contain the locally applied encoding for the targeted PHB.

All other ASes use this Active QoS Marking field to advertise their locally applied internal QoS encoding of the given class and technology at the interconnection point. The field value is right-aligned depending on the number of encoded bits. A cleared Marking field (all zero) signals that this traffic class experiences default traffic treatment within the transit AS forwarding technology.

4.3. Technology Type Enumeration

A small list of technologies is provided in the table below for the direct encoding of common technology types. The mapping of all virtual channel technologies into a single technology type value is for limiting the number of different communities in an UPDATE message. It is therefore a contribution to scalability.

Value	Technology Type
0x00	DiffServ enabled IP (DSCP encoding)
0x01	Ethernet using 802.1q priority tag
0x02	MPLS using E-LSP
0x03	Virtual Channel (VC) encoding using separate channels for QoS forwarding / one channel per class (e.g. ATM VCs, FR VCs, MPLS L-LSPs)
0x04	GMPLS - time slot encoding
0x05	GMPLS - lambda encoding
0x06	GMPLS - fibre encoding

Table 2

5. Community Usage

Providers MAY choose to process the QoS Marking communities and adopt the behaviour encoding and tunnel selection according to their local policy. Whether this MAY also lead to different IGP routing decisions or even effect BGP update filters is out of scope for the community definition.

Only the IP prefix originating AS is allowed to signal the QoS Marking communities and Sets. AS providers, which make use of this signalling mechanism MUST make sure, that only one external Class Set will be advertised for the AS. All advertised prefixes, which originate from that AS will be sent with the same QoS Marking community Set in the respective UPDATE message. Transit ASes MUST NOT modify or extend the QoS Marking community Set except for the update of each 'QoS Marking A' field contained in the community Set and the respective "P, R, I, A" flags. Prefixes with associated identical QoS Marking community Sets are to be advertised together in common UPDATE messages in relaying nodes.

Figure 4 shows an AS interconnection example with different Class Sets. It shows the case in AS 5 where different Class Sets are used internally and externally. The proposed QoS Class Set signalling will always use the external definitions within the UPDATE message QoS Marking communities. The example also shows, that IP prefixes, which originate in AS 5 and AS 3 can be advertised together with the same QoS Marking community Set as long as their Layer 2 encoding is identical.

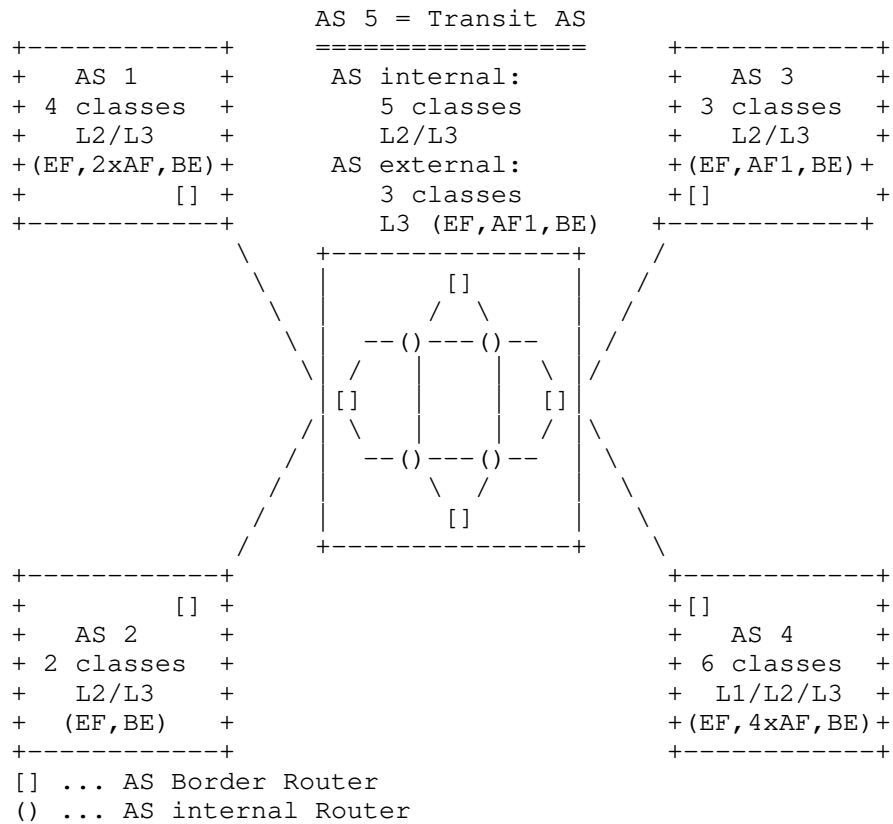


Figure 4

5.1. QoS Marking Example

See Appendix A for an example QoS Marking community Set.

5.2. AS Border Packet Forwarding

IP packet forwarding based on packet header QoS encoding might require remarking of packets in order to match AS internal policies and encodings of neighbouring ASes.

Identical QoS class sets and encodings between neighbouring ASes do not require any remarking. Different encodings will be matched on the outgoing traffic.

Outgoing traffic for a given IP prefix uses the 'QoS Marking A' information of the respective BGP UPDATE message QoS Marking community for adopted remarking of the forwarded packet.

If the 'I' flag is set for a given encoding, the outgoing traffic remarking SHOULD still be applied despite of the signalled lack of QoS Class forwarding support. This is particularly important, if the preserve flag 'P' is signalled together with the 'I' flag.

5.3. IP Prefix Aggregation

Several IP prefixes of different IP prefix originating ASes MAY be aggregated to a shorter IP prefix in transit ASes. If the original Class Sets of the aggregated prefixes are identical, the aggregate will use the same Set. In all other cases, the resulting IP prefix aggregate is handled the same as if the transit AS were the originating AS for this aggregated prefix. The transit AS provider MAY care for AS internal mechanisms, which map the signalled aggregate QoS Class Set to the different original Class Sets in the internal forwarding path.

In case of IP prefix aggregation with different QoS Class Sets, the 'Aggregation (A)' flag of each QoS Marking community within the Set MUST be set to '1'.

6. Confidentiality Considerations

The disclosure of confidential AS intrinsic information is of no concern since the signalled marking for QoS class encodings can be adopted prior to the UPDATE advertisement of the IP prefix originating AS. This way, a distinction between internal and external QoS Class Sets can be achieved. AS internal cross-layer marking adaptation and policy based update filtering allows for consistent QoS class support despite made up QoS Class Set and encoding information within UPDATE advertisements. In case of such policy hiding strategy, the required AS internal ingress and egress adaptation SHALL be done transparently without explicit "Active Marking" and 'R' flag signalling.

7. IANA Considerations

This document defines a new BGP Extended Community, which includes a "Technology Type" field. Section 4.3 enumerates a number of popular technologies. This list is expected to suffice for first implementations. However, future or currently uncovered technologies may arise, which will require an extended "Technology Type" enumeration list administered by IANA.

A new extended community QoS Marking community is defined, which has been assigned a Type value of 0x04 for a transitive and 0x44 for a non-transitive usage.

8. Security Considerations

This extension to BGP does not change the underlying security issues inherent in the existing BGP.

Malicious signalling behaviour of QoS Marking community advertising ASes can result in misguided neighbours about non existing or maliciously encoded Class Sets. Removal of QoS Marking community Sets leads to the current best effort interconnection, which is no stringent security concern.

The IP prefix originating AS MAY place a copy of its marking information into the Internet Routing Registry (IRR) for global reference.

The strongest threat is the advertisement of numerous very fine grained Class Sets, which could limit the scalability of this approach. However, neighbouring ASes are free to set the ignore flag of single communities or to stop processing the QoS Marking communities of a certain routing advertisement, once a self-set threshold has been crossed. By means of this self defence mechanism it should not be possible to crash neighbouring peers due to the excessive use of the new community.

9. References

9.1. Normative References

- [IANA_EC] IANA, "Border Gateway Protocol (BGP) Data Collection Standard Communities", June 2008, <<http://www.iana.org/assignments/bgp-extended-communities>>.
- [RFC1997] Chandra, R., Traina, P., and T. Li, "BGP Communities Attribute", RFC 1997, DOI 10.17487/RFC1997, August 1996, <<https://www.rfc-editor.org/info/rfc1997>>.
- [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>>.
- [RFC3140] Black, D., Brim, S., Carpenter, B., and F. Le Faucheur, "Per Hop Behavior Identification Codes", RFC 3140, DOI 10.17487/RFC3140, June 2001, <<https://www.rfc-editor.org/info/rfc3140>>.

- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, DOI 10.17487/RFC4271, January 2006, <<https://www.rfc-editor.org/info/rfc4271>>.
- [RFC4360] Sangli, S., Tappan, D., and Y. Rekhter, "BGP Extended Communities Attribute", RFC 4360, DOI 10.17487/RFC4360, February 2006, <<https://www.rfc-editor.org/info/rfc4360>>.
- [RFC5543] Ould-Brahim, H., Fedyk, D., and Y. Rekhter, "BGP Traffic Engineering Attribute", RFC 5543, DOI 10.17487/RFC5543, May 2009, <<https://www.rfc-editor.org/info/rfc5543>>.
- [RFC5575] Marques, P., Sheth, N., Raszuk, R., Greene, B., Mauch, J., and D. McPherson, "Dissemination of Flow Specification Rules", RFC 5575, DOI 10.17487/RFC5575, August 2009, <<https://www.rfc-editor.org/info/rfc5575>>.

9.2. Informative References

- [I-D.boucadair-qos-bgp-spec]
Boucadair, M., "QoS-Enhanced Border Gateway Protocol", draft-boucadair-qos-bgp-spec-01 (work in progress), July 2005.
- [I-D.jacquenet-bgp-qos]
Cristallo, G., "The BGP QOS_NLRI Attribute", draft-jacquenet-bgp-qos-00 (work in progress), February 2004.
- [I-D.knoll-idr-cos-interconnect]
Knoll, T., "BGP Class of Service Interconnection", draft-knoll-idr-cos-interconnect-22 (work in progress), May 2019.
- [I-D.liang-bgp-qos]
Benmohamed, L., "QoS Enhancements to BGP in Support of Multiple Classes of Service", draft-liang-bgp-qos-00 (work in progress), June 2006.
- [I-D.zhang-idr-bgp-extcommunity-qos]
Zhang, Z., "ExtCommunity map and carry TOS value of IP header", draft-zhang-idr-bgp-extcommunity-qos-00 (work in progress), November 2005.
- [MIT_CFP] Amante, S., Bitar, N., Bjorkman, N., and others, "Inter-provider Quality of Service - White paper draft 1.1", November 2006, <<http://cfp.mit.edu/docs/interprovider-qos-nov2006.pdf>>.

Appendix A. QoS Marking Example

The example AS is advertising several IP prefixes, which experience equal QoS treatment from AS internal networks. The IP packet forwarding policy within this originating AS defines e.g. 3 traffic classes for IP traffic (DSCP1, DSCP2 and DSCP3). These three classes are also consistently taken care of within a TC bit supporting MPLS tunnel forwarding. The BGP UPDATE message for the announced IP prefixes will contain the following QoS Marking community Set together with the IP prefix NLRI.

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	2	3	4	5	6	7	8	9
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

The class set as well as the example encodings are arbitrarily chosen.

Figure 5

Author's Address

Thomas Martin Knoll

Email: thomas.m.knoll@gmail.com