

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: 1 January 2023

K. Vairavakkalai, Ed.
N. Venkataraman
B. Rajagopalan
Juniper Networks, Inc.
G. Mishra
Verizon Communications Inc.
M. Khaddam
Cox Communications Inc.
X. Xu
Capitalonline.
R. Szarecki
Google.
J. Gowda
Extreme Networks
C. Yadlapalli
I. Means
AT&T
30 June 2022

BGP Classful Transport Planes
draft-kaliraj-idr-bgp-classful-transport-planes-17

Abstract

This document specifies a mechanism, referred to as "Intent Driven Service Mapping" to express association of overlay routes with underlay routes satisfying a certain SLA using BGP. The document describes a framework for classifying underlay routes into transport classes and mapping service routes to specific transport class.

The "Transport class" construct maps to a desired SLA and can be used to realize the "Topology Slice" in 5G Network slicing architecture.

This document specifies BGP protocol procedures that enable dissemination of such service mapping information that may span multiple cooperating administrative domains. These domains may be administered by the same provider or by closely co-ordinating provider networks.

A new BGP transport layer address family (SAFI 76) is defined for this purpose that uses RFC-4364 technology and follows RFC-8277 NLRI encoding. This new address family is called "BGP Classful Transport", aka BGP CT.

BGP CT makes it possible to advertise multiple tunnels to the same destination address, thus avoiding need of multiple loopbacks on the egress node.

It carries transport prefixes across tunnel domain boundaries (e.g. in Inter-AS Option-C networks), which is parallel to BGP LU (SAFI 4). It disseminates "Transport class" information for the transport prefixes across the participating domains, which is not possible with BGP LU. This makes the end-to-end network a "Transport Class" aware tunneled network.

Though BGP CT family is used only in the option-C inter-AS networks, the Service Mapping procedures described in this document apply in the same manner to Intra-AS service end points as well as Inter-AS option-A, option-B and option-C variations.

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 1 January 2023.

Copyright Notice

Copyright (c) 2022 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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	4
2. Terminology	6
3. Transport Class	8
4. "Transport Class" Route Target Extended Community	8
5. Transport Route Database	10
6. Nexthop Resolution Scheme	10
7. BGP Classful Transport Family NLRI	11
8. Use of Route Distinguisher	13
9. Comparison with other families using RFC-8277 encoding	13
10. Protocol Procedures	14
10.1. Preparing the network to deploy Classful Transport planes	14
10.2. Origination of Classful Transport route	15
10.3. Ingress node receiving Classful Transport route	15
10.4. Border node readvertising Classful Transport route with nexthop self	15
10.5. Border node receiving Classful Transport route on EBGP	16
10.6. Avoiding path-hiding through Route Reflectors	16
10.7. Avoiding loop between Route Reflectors in forwarding path	16
10.8. Ingress node receiving service route with Mapping Community	17
10.9. Coordinating between domains using different community namespaces	17
11. Flowspec Redirect to IP	18
12. BGP CT Egress TE	18
13. Interaction with BGP attributes specifying nexthop address and color	18
14. Scaling considerations	19
14.1. Avoiding unintended spread of BGP CT routes across domains	19
14.2. Constrained distribution of PNHs to SNs (On Demand Nexthop)	19
14.3. Limiting scope of visibility of PE loopback as PNHs	20
15. OAM considerations	21
16. Applicability to Network Slicing	22
17. SRv6 support	22
18. Illustration of procedures with example topology	23
18.1. Topology	23
18.2. Service Layer route exchange	24
18.3. Transport Layer route propagation	25
18.4. Data plane view	28
18.4.1. Steady state	28
18.4.2. Local repair of primary path	28

18.4.3. Absorbing failure of primary path. Fallback to best-effort tunnels.	29
19. IANA Considerations	29
19.1. New BGP SAFI	29
19.2. New Format for BGP Extended Community	30
19.2.1. Existing registries to be modified	30
19.2.2. New registries to be created	31
19.3. MPLS OAM code points	32
20. Security Considerations	32
21. Contributors	32
22. Acknowledgements	32
23. Normative References	32
Authors' Addresses	35

1. Introduction

The mechanisms defined in this document enable brownfield networks deployed using existing technologies like RSVP-TE and greenfield networks that use technologies like SPRING achieve 'Intent Driven Service Mapping'.

To facilitate this, the tunnels in a network can be grouped by the purpose they serve into a "Transport Class". These tunnels could be created using any signaling protocol including but not limited to LDP, RSVP-TE, BGP LU or SPRING. The tunnels could also use native IP or IPv6 as long as they can carry MPLS payload. Tunnels may exist between different pair of end points. Multiple tunnels may exist between the same pair of end points.

Thus, a Transport Class consists of tunnels created by various protocols that satisfy the properties of the class. For example, a "Gold" transport class may consist of tunnels that traverse the shortest path with fast re-route protection. A "Silver" transport class may hold tunnels that traverse shortest paths without protection. A "To NbrAS Foo" transport class may hold tunnels that exit to neighboring AS Foo and so on.

The extensions specified in this document can be used to create a BGP transport tunnel that potentially spans domains while preserving its Transport Class. Examples of domain are Autonomous System (AS) or IGP area. Within each domain, there is a second level underlay tunnel used by BGP to cross the domain. The second level underlay tunnels could be heterogeneous; each domain may use a different type of tunnel (e.g. MPLS, IP, GRE or SRv6) or use a different signaling protocol. A domain boundary is demarcated by a rewrite of BGP nexthop to 'self' while readvertising tunnel routes in BGP CT. Examples of domain boundary are inter-AS links and inter-region ABRs. The path uses MPLS label-switching when crossing domain boundaries and uses the native intra-AS tunnel of the desired transport class when traversing within a domain.

Overlay routes carry sufficient indication of the desired Transport Classes in the form of a BGP community called the "Mapping community". The "route resolution" procedure on the ingress node selects an appropriate tunnel whose destination matches (LPM) the nexthop of the overlay route belonging to the corresponding Transport Class. If the overlay route is carried in BGP, the protocol nexthop (or PNH) is carried as an attribute of the route.

The PNH of the overlay route is also referred to as "Service Endpoint" (SEP). The SEP may exist in the same domain as the service ingress node or lie in a different domain, which is adjacent or non-adjacent. In the former case, reachability to the SEP is provided by an intra-domain tunneling protocol and in the latter case, reachability to the SEP is via BGP transport families (e.g. SAFI 4 or 76).

In this architecture, the intra-domain transport protocols (e.g. RSVP-TE, SRTE) are also "Transport Class aware". They publish ingress routes in the Transport Route Database associated with the Transport Class at the tunnel ingress node. These routes are used to resolve BGP routes including BGP CT which may be further readadvertised to adjacent domains to extend this tunnel. How exactly the transport protocols are made transport class aware is outside the scope of this document.

This document describes mechanisms to:

- Model a "Transport Class" as a "Transport Route Database" on a router and to collect tunnel ingress routes of a certain class.

- Enable service routes to resolve over an intended Transport Class by virtue of carrying the appropriate "Mapping Community", which results in using the corresponding Transport Route Database for finding nexthop reachability.

Publish tunnel ingress routes in a Transport Route Database via BGP without any path hiding using BGP VPN and Add-path procedures, such that overlay routes in the receiving domains can also resolve over tunnels of the associated Transport Class.

Provide a way for cooperating domains to reconcile any differences in extended community namespaces and interoperate between different transport signaling protocols in each domain.

In this document we focus mainly on MPLS as the intra-domain transport tunnel forwarding technology, but the mechanisms described here would work in similar manner for non-MPLS (e.g. IP, GRE, UDP or SRv6) transport tunnel forwarding technologies too.

This document assumes MPLS forwarding as the defacto standard when crossing domain boundaries. However mechanisms specified in this document can also support different forwarding technologies (e.g. SRv6). Section 17 (SRv6 support) in this document describes the application of BGP CT over SRv6 data plane.

The document Seamless Segment Routing [Seamless-SR] describes various use cases and applications of procedures described in this document.

2. Terminology

LSP: Label Switched Path.

TE : Traffic Engineering.

SN : Service Node. A router that sends or receives BGP Service routes (e.g. SAFI 1, 128) with self as nexthop.

BN : Border Node. A router that sends or receives BGP Transport routes (e.g. SAFI 4, 76) with self as nexthop.

TN : Transport Node, P-router.

BGP-VPN : VPNs built using RFC4364 mechanisms.

RT : Route-Target extended community.

RD : Route-Distinguisher.

VRF: Virtual Router Forwarding Table.

CsC: Carrier serving Carrier VPN.

PNH : Protocol-Nexthop address carried in a BGP Update message.

EP : End point, a loopback address in the network.

SEP : Service End point, the PNH of a Service route.

LPM : Longest Prefix Match.

SLA: Service Level Agreement.

EPE: Egress Peer Engineering.

Service Family : BGP address family used for advertising routes for "data traffic" as opposed to tunnels (e.g. SAFI 1 or 128).

Transport Family : BGP address family used for advertising tunnels, which are in turn used by service routes for resolution (e.g. SAFI 4 or 76).

Transport Tunnel : A tunnel over which a service may place traffic (e.g. GRE, UDP, LDP, RSVP-TE or SPRING).

Tunnel Ingress Route: Route to Tunnel Destination/Endpoint installed at the headend (ingress) of the tunnel by the tunneling protocol.

Tunnel Domain : A domain of the network containing SNs and BNs under a single administrative control that has tunnels between them. An end-to-end tunnel spanning several adjacent tunnel domains can be created by "stitching" them together using labels.

Transport Class : A group of transport tunnels offering the same SLA.

Transport Class RT : A Route-Target extended community used to identify a specific Transport Class.

Transport Route Database : At the SN and BN, a Transport Class has an associated Transport Route Database that collects its tunnel ingress routes.

Transport Plane : An end to end plane comprising of transport tunnels belonging to same Transport Class. Tunnels of same Transport Class are stitched together by BGP CT route readvertisements with nexthop self to enable Label-Swap forwarding across domain boundaries.

Mapping Community : BGP Community/Extended-community on a service route that maps it to resolve over a Transport Class.

3. Transport Class

A Transport Class is defined as a set of transport tunnels that share the same SLA. It is encoded as the Transport Class RT, which is a new Route-Target extended community.

A Transport Class is configured at SN and BN with RD and Route Target attributes. Creation of a Transport Class instantiates its corresponding Transport Route Database.

The operator may configure an SN/BN to classify a tunnel into an appropriate Transport Class, which causes the tunnel's ingress route to be installed in the corresponding Transport Route Database. These routes are used to resolve BGP routes including BGP CT which may be further readvertised to adjacent domains to extend this tunnel.

Alternatively, a router receiving the transport routes in BGP with appropriate signaling information can associate those ingress routes to the appropriate Transport Class. E.g. for Classful Transport family (SAFI 76) routes, the Transport Class RT indicates the Transport Class. For BGP LU family (SAFI 4) routes, import processing based on Communities or inter-AS source-peer may be used to place the route in the desired Transport Class.

When the ingress route is received via SRTE [SRTE] with "Color:Endpoint" as the NLRI that encodes the Transport Class as an integer 'Color', the 'Color' is mapped to a Transport Class during import processing. SRTE ingress route for 'Endpoint' is installed in the corresponding Transport Route Database. The SRTE tunnel will be extended by a BGP CT advertisement with NLRI 'RD:Endpoint', Transport Class RT and a new label. The MPLS swap route thus installed for the new label will pop the label and deliver decapsulated traffic into the path determined by SRTE route.

RFC8664 [RFC8664] extends PCEP to carry SRTE Color. This color association learnt from PCEP is also mapped to a Transport Class thus associating the PCEP signaled SRTE LSP with the desired Transport Class.

Similarly, PCEP-RSVP-COLOR [PCEP-RSVP-COLOR] extends PCEP to carry RSVP Color. This color association learnt from PCEP is also mapped to a Transport Class thus associating the PCEP signaled RSVP-TE LSP with the desired Transport Class.

4. "Transport Class" Route Target Extended Community

This document defines a new type of Route Target, called "Transport Class" Route Target Extended Community.

"Transport Class" Route Target extended community is a transitive extended community EXT-COMM [RFC4360] of extended-type, with a new Format (Type high = 0xa) and SubType as 0x2 (Route Target).

This new Route Target Format has the following encoding:

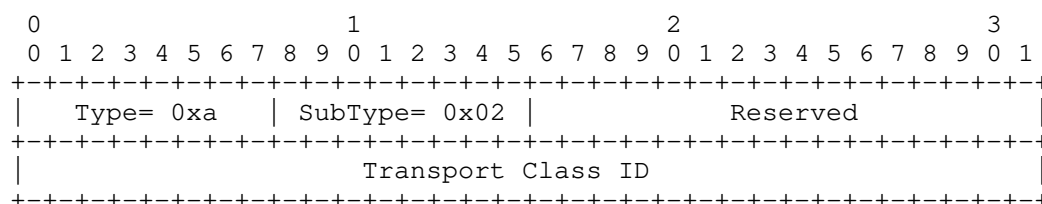


Fig 1: "Transport Class" Route Target Extended Community

Type: 1 octet

Type field contains value 0xa.

SubType: 1 octet

Subtype field contain 0x2. This indicates 'Route Target'.

Transport Class ID: 4 octets

The least significant 32-bits of the value field contain the "Transport Class" identifier, which is a 32-bit integer.

The remaining 2 octets after SubType field are Reserved. They MUST be set to zero on transmission, SHOULD be ignored on reception and left unaltered.

The "Transport class" Route Target Extended community follows the mechanisms for VPN route import/export as specified in BGP-VPN [RFC4364] and follows the Constrained Route Distribution mechanisms as specified in Route Target Constraints [RFC4684]

A BGP speaker that implements RT Constraint Route Target Constraints [RFC4684] MUST apply the RT Constraint procedures to the "Transport class" Route Target Extended community as well.

The Transport Class Route Target Extended community is carried on Classful Transport family routes and allows associating them with appropriate Transport Route Databases at receiving BGP speakers.

Use of the Transport Class Route Target Extended community with a new Type code avoids conflicts with any VPN Route Target assignments already in use for service families.

5. Transport Route Database

A Transport Route Database is a logical collection of transport routes pertaining to the same Transport Class. Tunnel endpoint addresses in this database belong to the "Provider Namespace".

Overlay routes that want to use a specific Transport Class confine the scope of nexthop resolution to the set of routes contained in the corresponding Transport Route Database.

The Transport Route Database can be realized as a "Routing Table" referred in Section 9.1.2.1 of RFC4271 (<https://www.rfc-editor.org/rfc/rfc4271#section-9.1.2.1>) which is a control plane only database. However, an implementation may choose a different methodology to realize this logical construct in such a way that it supports the procedures defined in this document.

SN or BN originate routes for 'Classful Transport' address family from the Transport Route Database. These routes have NLRI "RD:Endpoint", Transport Class RT and an MPLS label. 'Classful Transport' family routes received with Transport Class RT are imported into its corresponding Transport Route Database.

6. Nexthop Resolution Scheme

An implementation may provide an option for the service route to resolve over less preferred Transport Classes, should the resolution over preferred or "primary" Transport Class fail.

To accomplish this, the set of service routes may be associated with a user-configured "Resolution Scheme" that consists of the primary Transport Class and an optional ordered list of fallback Transport Classes.

A community called as "Mapping Community" is configured for a "resolution scheme". A Mapping Community maps to exactly one Resolution Scheme. A Resolution Scheme comprises of one primary transport class and optionally, one or more fallback transport classes.

A BGP route is associated with a resolution scheme during import processing. The first community on the route that matches a Mapping Community of a locally configured Resolution Scheme is considered the effective Mapping Community for the route. The Resolution Scheme

thus found is used when resolving the route's PNH. If a route contains more than one Mapping Community, it indicates that the route considers these multiple Mapping Communities as equivalent. So, the first community that maps to a Resolution Scheme is chosen.

A transport route received in BGP Classful Transport family SHOULD use a Resolution Scheme that contains the primary Transport Class without any fallback to best effort tunnels. The primary Transport Class is identified by the Transport Class RT carried on the route. Thus, Transport Class RT serves as the Mapping Community for BGP CT routes.

A service route received in a BGP service family MAY map to a Resolution Scheme that contains the primary Transport Class identified by the Mapping Community on the route and a fallback to best effort Transport Class. The primary Transport Class is identified by the Mapping Community carried on the route. For e.g. the Extended Color community may serve as the Mapping Community for service routes. Color:0:<n> MAY map to a Resolution Scheme that has primary Transport Class <n> and a fallback to best-effort Transport Class.

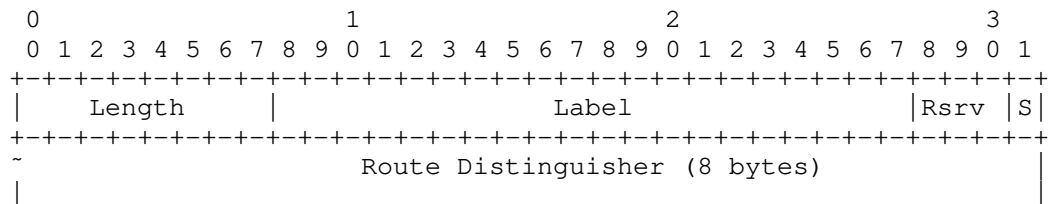
7. BGP Classful Transport Family NLRI

The Classful Transport (CT) family will use the existing AFI of IPv4 or IPv6 and a new SAFI 76 "Classful Transport" that will apply to both IPv4 and IPv6 AFIs. These AFI, SAFI pair of values MUST be negotiated in Multiprotocol Extensions capability described in [RFC4760] to be able to send and receive BGP CT routes.

The "Classful Transport" SAFI NLRI itself is encoded as specified in <https://tools.ietf.org/html/rfc8277#section-2> [RFC8277].

When AFI/SAFI is 1/76, the Classful Transport NLRI Prefix consists of an 8-byte RD followed by an IPv4 prefix. When AFI/SAFI is 2/76, the Classful Transport NLRI Prefix consists of an 8-byte RD followed by an IPv6 prefix.

For better readability, the following figure illustrates a BGP Classful Transport family NLRI when single Label is advertised:



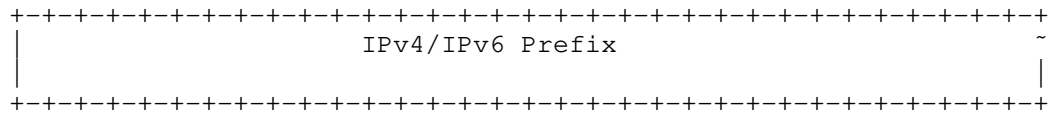


Fig 2: SAFI 76 "Classful Transport" NLRI

Length: 1 octet

The Length field consists of a single octet. It specifies the length in bits of the remainder of the NLRI field.

Note that the length will always be the sum of 20 (number of bits in Label field), plus 3 (number of bits in Rsrv field), plus 1 (number of bits in S field), plus the length in bits of the Prefix (RD:IP prefix).

In an MP_REACH_NLRI attribute whose SAFI is 76, the Prefix (RD + IP prefix) will be 96 bits or less if the AFI is 1 and will be 192 bits or less if the AFI is 2.

As specified in [RFC4760], the actual length of the NLRI field will be the number of bits specified in the Length field, rounded up to the nearest integral number of octets.

Label:

The Label field is a 20-bit field containing an MPLS label value (see [RFC3032]).

Rsrv:

This 3-bit field SHOULD be set to zero on transmission and MUST be ignored on reception.

S:

When single label is advertised, this 1-bit field MUST be set to one on transmission and MUST be ignored on reception.

Route Distinguisher:

8 byte RD as defined in [RFC4364 Sec 4.2].

IPv4/IPv6 Prefix:

IPv4 prefix, if AFI/SAFI 1/76.
IPv6 prefix, if AFI/SAFI 2/76.

Attributes on a Classful Transport route include the Transport Class Route-Target extended community, which is used to associate the route with the correct Transport Route Databases on SNs and BNs in the network.

SAFI 76 routes can be sent with either IPv4 or IPv6 nexthop. The type of nexthop is inferred from the length of the nexthop.

When the length of Next Hop Address field is 24 (or 48) the nexthop address is of type VPN-IPv6 with 8-octet RD set to zero (potentially followed by the link-local VPN-IPv6 address of the next hop with an 8-octet RD set to zero).

When the length of the Next Hop Address field is 12 the nexthop address is of type VPN-IPv4 with 8-octet RD set to zero.

8. Use of Route Distinguisher

RD aids in troubleshooting a BGP CT network by uniquely identifying the originator of a route across a multi domain network.

Use of RD also allows the option for signaling forwarding diversity within the same Transport Class. The same Egress PE can advertise multiple BGP CT routes for an EP belonging to the same Transport Class.

E.g. multiple RDx:EP1 prefixes can be advertised for an EP1 to different set of BGP peers in order to collect traffic statistics for them. In absense of RD, duplicated Transport Class/Color values will be needed in the transport network to achieve such use cases.

9. Comparison with other families using RFC-8277 encoding

SAFI 128 (Inet-VPN) is an RFC8277 encoded family that carries service prefixes in the NLRI, where the prefixes come from the customer namespaces and are contextualized into separate user virtual service RIBs called VRFs using RFC4364 procedures.

SAFI 4 (BGP LU) is an RFC8277 encoded family that carries transport prefixes in the NLRI, where the prefixes come from the provider namespace.

SAFI 76 (Classful Transport) is an RFC8277 encoded family that carries transport prefixes in the NLRI, where the prefixes come from the provider namespace and are contextualized into separate Transport Route Databases using RFC4364 procedures.

It is worth noting that SAFI 128 has been used to carry transport prefixes in "L3VPN Inter-AS Carrier's carrier" scenario, where BGP LU/LDP prefixes in CsC VRF are advertised in SAFI 128 towards the remote-end client carrier.

In this document a new AFI/SAFI is used instead of reusing SAFI 128 to carry these transport routes because it is operationally advantageous to segregate transport and service prefixes into separate address families. E.g. It allows to safely enable "per-prefix" label allocation scheme for Classful Transport prefixes without affecting SAFI 128 service prefixes which may have huge scale. The "per prefix" label allocation scheme keeps the routing churn local during topology changes.

A new family also facilitates having a different readvertisement path of the transport family routes in a network than the service route readvertisement path. Service routes (Inet-VPN) are exchanged over an EBGP multihop session between Autonomous systems with nexthop unchanged; whereas Classful Transport routes are readvertised over EBGP single hop sessions with "nexthop self" rewrite over inter-AS links.

The Classful Transport family is similar in vein to BGP LU, in that it carries transport prefixes. The only difference is that it also carries in Route Target, an indication of which Transport Class the transport prefix belongs to and uses RD to disambiguate multiple instances of the same transport prefix in a BGP Update.

10. Protocol Procedures

This section summarizes the procedures followed by various nodes speaking Classful Transport family.

10.1. Preparing the network to deploy Classful Transport planes

Operator decides on the Transport Classes that exist in the network and allocates a Transport Class Route Target to identify each Transport Class.

Operator configures Transport Classes on the SNs and BNs in the network with Transport Class Route Targets and unique Route-Distinguishers.

Implementations MAY provide automatic generation and assignment of RD, RT values; they MAY also provide a way to manually override the automatic mechanism in order to deal with any conflicts that may arise with existing RD, RT values in different network domains participating in the deployment.

10.2. Origination of Classful Transport route

At the ingress node of the tunnel's home domain, the tunneling protocols install tunnel ingress routes in the Transport Route Database associated with the Transport Class the tunnel belongs to.

The egress node of the tunnel i.e. the tunnel endpoint originates the BGP Classful Transport route with NLRI containing RD:TunnelEndpoint, Transport Class RT and PNH TunnelEndpoint, which will resolve over the tunnel route in Transport Route Database at the ingress node. When the tunnel is up, the Classful Transport BGP route will become usable and get re-advertised.

Alternatively, the ingress node may advertise this tunnel destination into BGP as a Classful Transport family route with NLRI RD:TunnelEndpoint, attaching a 'Transport Class' Route Target that identifies the Transport Class. This BGP CT route is advertised to EBGp peers and IBGP peers in neighboring domains. This route SHOULD NOT be advertised to the IBGP core that contains the tunnel.

Unique RD SHOULD be used by the originator of a Classful Transport route to disambiguate the multiple BGP advertisements for a transport end point.

10.3. Ingress node receiving Classful Transport route

On receiving a BGP Classful Transport route with a PNH that is not directly connected (e.g. an IBGP-route), a Mapping Community on the route (the Transport Class RT) indicates which Transport Class this route maps to. The routes in the associated Transport Route Database are used to resolve the received PNH. If there does not exist a route in the Transport Route Database matching the PNH, the Classful Transport route is considered unusable and MUST NOT be advertised further.

10.4. Border node readvertising Classful Transport route with nexthop self

The BN allocates an MPLS label to advertise upstream in Classful Transport NLRI. The BN also installs an MPLS route for that label that swaps the incoming label with a label received from the downstream BGP speaker or pops the incoming label. It then pushes received traffic to the transport tunnel or direct interface that the Classful Transport route's PNH resolved over.

The label SHOULD be allocated with "per-prefix" label allocation semantics. RD is stripped from the BGP CT NLRI prefix when a BGP CT route is added to a Transport Route Database. The IP prefix in the Transport Route Database context (Transport-Class, IP-prefix) is used as the key to do per-prefix label allocation. This helps in avoiding BGP CT route churn through out the CT network when a failure happens in a domain. The failure is not propagated further than the BN closest to the failure.

The value of advertised MPLS label is locally significant, and is dynamic by default. The BN may provide option to allocate a value from a statically carved out range. This can be achieved using locally configured export policy, or via mechanisms described in BGP Prefix-SID [RFC8669].

10.5. Border node receiving Classful Transport route on EBGp

If the route is received with PNH that is known to be directly connected (e.g. EBGp single-hop peering address), the directly connected interface is checked for MPLS forwarding capability. No other nexthop resolution process is performed as the inter-AS link can be used for any Transport Class.

If the inter-AS links should honor Transport Class, then the BN SHOULD follow procedures of an Ingress node described above and perform nexthop resolution process. The interface routes SHOULD be installed in the Transport Route Database belonging to the associated Transport Class.

10.6. Avoiding path-hiding through Route Reflectors

When multiple BNs exist such that they advertise a RD:EP prefix to RRs, the RRs may hide all but one of the BNs, unless ADDPATH [RFC7911] is used for the Classful Transport family. This is similar to L3VPN option-B scenarios. Hence ADDPATH SHOULD be used for Classful Transport family, to avoid path-hiding through RRs.

10.7. Avoiding loop between Route Reflectors in forwarding path

Pair of redundant ABRs, each acting as an RR with nexthop self may choose each other as best path instead of the upstream ASBR, causing a traffic forwarding loop.

Implementations SHOULD provide a way to alter the tie-breaking rule specified in BGP RR [RFC4456] to tie-break on CLUSTER_LIST step before ROUTER-ID step, when performing path selection for BGP CT routes. RFC4456 considers pure RR which is not in forwarding path. When RR is in forwarding path and reflects routes with

nexthop self as is the case for ABR BNs in a BGP transport network, this rule may cause loops. This document suggests the following modification to the BGP Decision Process Tie Breaking rules (Sect. 9.1.2.2, [RFC4271]) when doing path selection for BGP CT family routes:

The following rule SHOULD be inserted between Steps e) and f): a BGP Speaker SHOULD prefer a route with the shorter CLUSTER_LIST length. The CLUSTER_LIST length is zero if a route does not carry the CLUSTER_LIST attribute.

Some deployment considerations can also help in avoiding this problem:

- IGP metric should be assigned such that "ABR to redundant ABR" cost is inferior than "ABR to upstream ASBR" cost.
- Tunnels belonging to non best effort Transport Classes SHOULD NOT be provisioned between ABRs. This will ensure that the route received from an ABR with nexthop self will not be usable at a redundant ABR.

This avoids possibility of such loops altogether.

10.8. Ingress node receiving service route with Mapping Community

Service routes received with Mapping Community resolve using Transport Route Databases determined by the Resolution Scheme. If the resolution process does not find a Tunnel Ingress Route in any of the Transport Route Databases, the service route MUST be considered unusable for forwarding purpose and be withdrawn.

10.9. Coordinating between domains using different community namespaces

Cooperating option-C domains may sometimes not agree on RT, RD, Mapping-community or Transport Route Target values because of differences in community namespaces (e.g. during network mergers or renumbering for expansion). Such deployments may deploy mechanisms to map and rewrite the Route Target values on domain boundaries, using per ASBR import policies. This is no different than any other BGP VPN family. Mechanisms used in inter-AS VPN deployments may be used with the Classful Transport family also.

The Resolution Schemes SHOULD allow association with multiple Mapping Communities. This helps with renumbering, network mergers or transitions.

Deploying unique RDs is strongly RECOMMENDED because it helps in troubleshooting by uniquely identifying originator of a route and avoids path-hiding.

This document defines a new format of Route-Target extended-community to carry Transport Class, this avoids collision with regular Route Target namespace used by service routes.

11. Flowspec Redirect to IP

Flowspec routes using Redirect to IP nexthop is described in BGP Flow-Spec Redirect to IP Action [FLOWSPEC-REDIR-IP]

Such Flowspec BGP routes with Redirect to IP nexthop can be attached with a Mapping Community (e.g. Color:0:100), which allows redirecting the flow traffic over a tunnel to the IP nexthop satisfying the desired SLA (e.g. Transport Class color 100).

Flowspec BGP family acts as just another service that can make use of BGP CT architecture to achieve Flow based forwarding with SLAs.

12. BGP CT Egress TE

Mechanisms described in BGP LU EPE [BGP-LU-EPE] also applies to BGP CT family.

The Peer/32 or Peer/128 EPE route MAY be originated in BGP CT family with appropriate Mapping Community (e.g. transport-target:0:100), thus allowing an EPE path to the peer that satisfies the desired SLA.

13. Interaction with BGP attributes specifying nexthop address and color

The Tunnel Encapsulation Attribute described in RFC9012 [RFC9012] can be used to request a specific type of tunnel encapsulation. Usage of this attribute may apply to BGP service routes or transport routes, including BGP Classful Transport family routes.

Mechanisms described in BGP MultiNexthop Attribute [MULTI-NH-ATTR] allow a BGP route to carry multiple nexthop addresses. It also allows specifying 'Transport Class ID' as a qualifier for each Nexthop address.

It should be noted that in such cases "Transport Class/Color" can exist in multiple places on the same route, and a precedence order needs to be established to determine which Transport class the route's nexthop should resolve over. This document suggests the following order of precedence, more preferred first:

Transport Class ID SubTLV, in MultiNexthop Attribute.

Color SubTLV, in Tunnel Encapsulation Attribute.

Transport Target Extended community, on BGP CT route.

Color Extended community, on BGP service route.

The above precedence order follows more specific scoping of Color to less specific scoping.

Transport Class ID specified for Nexthop-Leg subTLV in a MultiNextHop attribute is more specific indication of Color than Color subTLV in a TEA, which inturn is more specific than Mapping Community (Transport Target) on a BGP CT transport route, which is inturn more specific than a Service route scoped Mapping Community (Color Extended community).

14. Scaling considerations

14.1. Avoiding unintended spread of BGP CT routes across domains

RFC8212 [RFC8212] suggests BGP speakers require explicit configuration of both BGP Import and Export Policies in order to receive or send routes over EBGP sessions.

It is recommended to follow this for BGP CT routes. It will prohibit unintended advertisement of transport routes throughout the BGP CT transport domain which may span across multiple AS domains. This will conserve usage of MPLS label and nexthop resources in the network. An ASBR of a domain can be provisioned to allow routes with only the Transport Route Targets that are required by SNs in the domain.

14.2. Constrained distribution of PNHs to SNs (On Demand Nexthop)

This section describes how the number of Protocol Nexthops advertised to a SN or BN can be constrained using BGP Classful Transport and Route Target Constraints [RFC4684].

An egress SN MAY advertise BGP CT route for RD:eSN with two Route Targets: transport-target:0:<TC> and a RT carrying <eSN>:<TC>. Where TC is the Transport Class identifier, and eSN is the IP-address used by SN as BGP nexthop in its service route advertisements.

The transport-target:0:<TC> is the new type of route target (Transport Class RT) defined in this document. It is carried in BGP extended community attribute (BGP attribute code 16).

The RT carrying <eSN>:<TC> MAY be an IP-address specific regular RT (BGP attribute code 16), IPv6-address specific RT (BGP attribute code 25), or a Wide-communities based RT (BGP attribute code 34) as described in Route Target Constrain Extension [RTC-Ext]. This document recommends using Wide-communities based RT for the same.

An ingress SN MAY import BGP CT routes with Route Target carrying <eSN>:<TC>. The ingress SN MAY learn the eSN values either by configuration, or it MAY discover them from the BGP nexthop field in the BGP VPN service routes received from eSN. A BGP ingress SN receiving a BGP service route with nexthop of eSN SHOULD generate a RTC/Extended-RTC route for Route Target prefix <Origin ASN>:<eSN>/[80|176] in order to learn BGP CT transport routes to reach eSN. This allows constrained distribution of the transport routes to the PNHs actually required by iSN.

When path of route propagation of BGP CT routes is same as the RTC routes, a BN would learn the RTC routes advertised by ingress SNs and propagate further. This will allow constraining distribution of BGP CT routes for a PNH to only the necessary BNs in the network, closer to the egress SN.

This mechanism provides "On Demand Nexthop" of BGP CT routes, which help with the scaling of MPLS forwarding state at SN and BN.

However, the amount of state carried in RTC family may become proportional to number of PNHs in the network. To strike a balance, the RTC route advertisements for <Origin ASN>:<eSN>/[80|176] MAY be confined to the BNs in home region of ingress-SN, or the BNs of a super core.

Such a BN in the core of the network SHOULD import BGP CT routes with Transport-Target:0:<TC> and generate a RTC route for <Origin ASN>:0:<TC>/96, while not propagating the more specific RTC requests for specific PNHs. This will let the BN learn transport routes to all eSN nodes. But confine their propagation to ingress-SNs.

14.3. Limiting scope of visibility of PE loopback as PNHs

It may be even more desirable to limit the number of PNHs that are globally visible in the network. This is possible using mechanism described in MPLS Namespaces [MPLS-NAMESPACES]

Such that advertisement of PE loopback addresses as next-hop in BGP service routes is confined to the region they belong to. An anycast IP-address called "Context Protocol Nexthop Address" (CPNH) abstracts the SNs in a region from other regions in the network, swapping the SN scoped service label with a CPNH scoped private namespace label.

This provides much greater advantage in terms of scaling and convergence. Changes to implement this feature are required only on the local region's BNs and RRs.

15. OAM considerations

Standard MPLS OAM procedures specified in [RFC8029] also apply to BGP Classful Transport.

The 'Target FEC Stack' sub-TLV for IPv4 Classful Transport has a Sub-Type of [TBD], and a length of 13. The Value field consists of the RD advertised with the Classful Transport prefix, the IPv4 prefix (with trailing 0 bits to make 32 bits in all) and a prefix length encoded as follows:

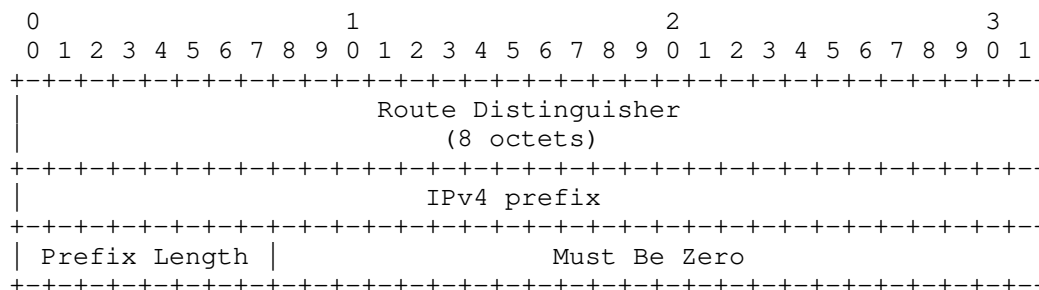


Figure 1: Classful Transport IPv4 FEC

The 'Target FEC Stack' sub-TLV for IPv6 Classful Transport has a Sub-Type of [TBD], and a length of 25. The Value field consists of the RD advertised with the Classful Transport prefix, the IPv6 prefix (with trailing 0 bits to make 128 bits in all) and a prefix length encoded as follows:

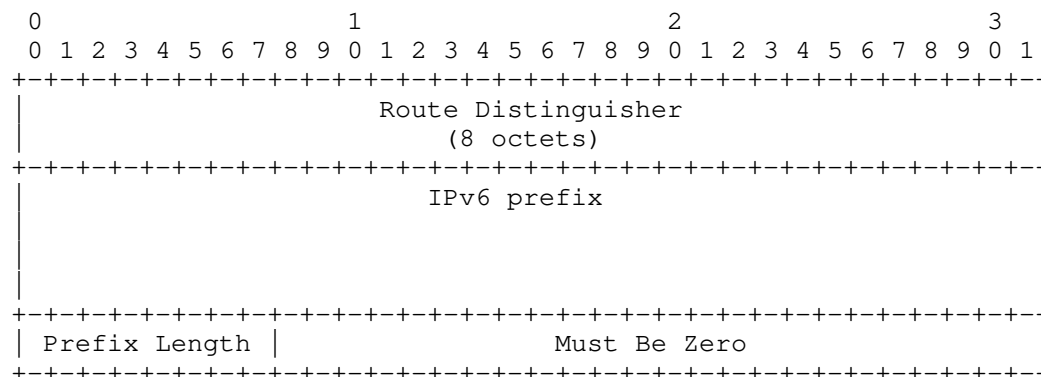


Figure 2: Classful Transport IPv6 FEC

16. Applicability to Network Slicing

In Network Slicing, the Transport Slice Controller (TSC) sets up the Topology (e.g. RSVP-TE, SR-TE tunnels with desired characteristics) and resources (e.g. polices/shapers) in a transport network to create a Transport Slice. The Transport Class construct described in this document represents the "Topology Slice" portion of this equation.

The TSC can use the Transport Class Identifier (Color value) to provision a transport tunnel in a specific Topology Slice.

Further, Network Slice Controller can use the Mapping Community on the service route to map traffic to the desired Transport Slice.

17. SRv6 support

This section describes how BGP CT may be used to set up inter domain tunnels of a certain Transport Class, when using Segment Routing over IPv6 (SRv6) data plane on the inter AS links or as an intra AS tunneling mechanism.

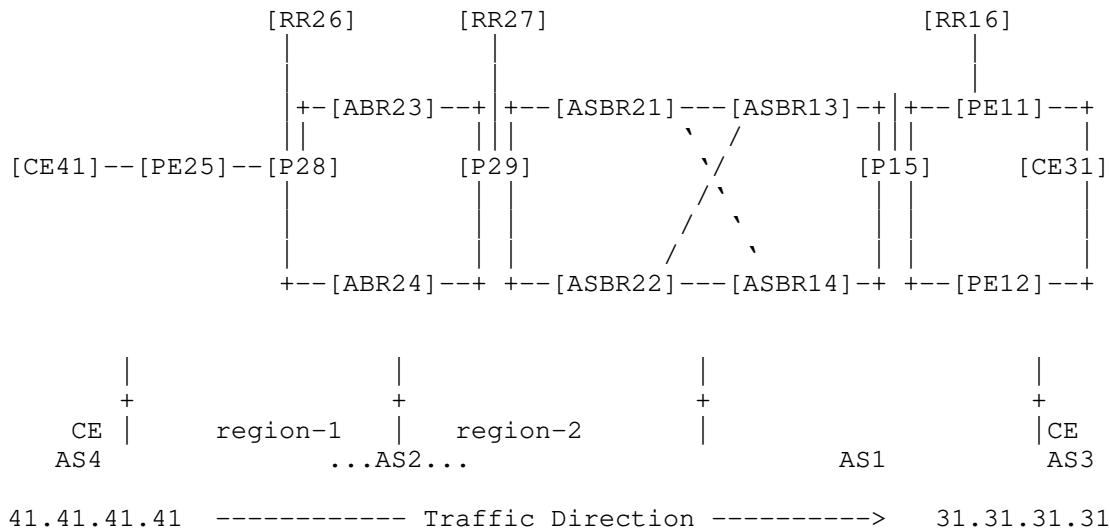
[RFC8986] specifies the SRv6 Endpoint behaviors (End.USD, End.BM, End.B6.Encaps and End.Replace, End.ReplaceB6, respectively). [SRV6-INTER-DOMAIN] specify the SRv6 Endpoint behaviors (END.REPLACE, END.REPLACEB6 and END.DB6). These are leveraged for BGP CT with SRv6 data plane.

The BGP Classful Transport route update for SRv6 MUST include the BGP Prefix-SID attribute along with SRv6 SID information as specified in [SRV6-SERVICES]. It may also include SRv6 SID structure for Transposition as specified in [SRV6-SERVICES]. It should be noted that prefixes carried in BGP CT family are transport layer end-points, e.g. PE loopback addresses. Thus the SRv6 SID carried in a BGP CT route is also a transport layer identifier.

This document extends the usage of "SRv6 label route tunnel" TLV to AFI=1/2 SAFI 76. "SRv6 label route tunnel" is the TLV of the BGP Prefix-SID Attribute as specified in [SRV6-MPLS-AGRWL].

18. Illustration of procedures with example topology

18.1. Topology



This example shows a provider network that comprises of two Autonomous systems, AS1, AS2. They are serving customers AS3, AS4 respectively. Traffic direction being described is CE41 to CE31. CE31 may request a specific SLA (e.g. Gold for this traffic), when traversing these provider networks.

AS2 is further divided into two regions. So, there are three tunnel domains in provider space. AS1 uses ISIS Flex-Algo intra-domain tunnels, whereas AS2 uses RSVP-TE intra-domain tunnels.

The network has two Transport classes: Gold with transport class id 100, Bronze with transport class id 200. These transport classes are provisioned at the PEs and the Border nodes (ABRs, ASBRs) in the network.

Following tunnels exist for Gold transport class.

PE25_to_ABR23_gold - RSVP-TE tunnel

PE25_to_ABR24_gold - RSVP-TE tunnel

ABR23_to_ASBR22_gold - RSVP-TE tunnel

ASBR13_to_PE11_gold - ISIS FlexAlgo tunnel

ASBR14_to_PE11_gold - ISIS FlexAlgo tunnel

Following tunnels exist for Bronze transport class.

PE25_to_ABR23_bronze - RSVP-TE tunnel

ABR23_to_ASBR21_bronze - RSVP-TE tunnel

ABR23_to_ASBR22_bronze - RSVP-TE tunnel

ABR24_to_ASBR21_bronze - RSVP-TE tunnel

ASBR13_to_PE12_bronze - ISIS FlexAlgo tunnel

ASBR14_to_PE11_bronze - ISIS FlexAlgo tunnel

These tunnels are either provisioned or auto-discovered to belong to transport class 100 or 200.

18.2. Service Layer route exchange

Service nodes PE11, PE12 negotiate service families (SAFI 1, 128) on the BGP session with RR16. Service helpers RR16, RR26 have multihop EBGp session to exchange service routes between the two AS. Similarly PE25 negotiates service families with RR26.

Forwarding happens using service routes at service nodes PE25, PE11, PE12 only. Routes received from CEs are not present in any other nodes' FIB in the network.

CE31 advertises a route for example prefix 31.31.31.31 with nexthop self to PE11, PE12. CE31 can attach a Mapping Community Color:0:100 on this route, to indicate its request for Gold SLA. Or, PE11 can attach the same using locally configured policies.

Consider CE31 is getting VPN service from PE11. The 31.31.31.31 route is readvertised in SAFI 128 by PE11 with nexthop self (1.1.1.1) and label V-L1, to RR16 with the Mapping Community Color:0:100 attached. This SAFI 128 route reaches PE25 via RR16, RR26 with the nexthop unchanged, as PE11 and label V-L1. Now PE25 can resolve the PNH 1.1.1.1 using transport routes received in BGP CT or BGP LU.

The IP FIB at PE25 VRF will have a route for 31.31.31.31 with a nexthop when resolved, that points to a Gold tunnel in ingress domain.

18.3. Transport Layer route propagation

Egress nodes PE11, PE12 negotiate BGP CT family with transport ASBRs ASBR13, ASBR14. These egress nodes originate BGP CT routes for tunnel endpoint addresses, that are advertised as nexthop in BGP service routes. In this example both PEs participate in transport classes Gold and Bronze. The protocol procedures are explained using Gold SLA plane and the Bronze SLA plane is used to highlight the path hiding aspects.

PE11 is provisioned with transport class 100, RD value 1.1.1.1:10 and a transport-target:0:100 for Gold tunnels. And a Transport class 200 with RD value 1.1.1.1:20, and transport route target 0:200 for Bronze tunnels. Similarly, PE12 is provisioned with transport class 100, RD value 1.1.1.2:10 and a transport-target:0:100 for Gold tunnels. And transport class 200, RD value 1.1.1.2:20 with transport-target:0:200 for Bronze tunnels

Similarly, these transport classes are also configured on ASBRs, ABRs and PEs with same Transport Route Target and unique RDs.

ASBR13 and ASBR14 negotiate BGP CT family with transport ASBRs ASBR21, ASBR22 in neighboring AS. They negotiate BGP CT family with RR27 in region 2, which reflects BGP CT routes to ABR23, ABR24. ABR23, ABR24 negotiate BGP CT family with Ingress node PE25 in region 1. BGP LU family is also negotiated on these sessions alongside BGP CT family. BGP LU carries "best effort" transport class routes, BGP CT carries gold, bronze transport class routes.

PE11 is provisioned to originate BGP CT route with Gold SLA to endpoint PE11. This route is sent with NLRI RD prefix 1.1.1.1:10:1.1.1.1, Label B-L0, nexthop 1.1.1.1 and a route target extended community transport-target:0:100. Label B-L0 can either be Implicit Null (Label 3) or a Ultimate Hop Pop (UHP) label.

This route is received by ASBR13 and it resolves over the tunnel ASBR13_to_PE11_gold. The route is then readvertised by ASBR13 in BGP CT family to ASBRs ASBR21, ASBR22 according to export policy. This route is sent with same NLRI RD prefix 1.1.1.1:10:1.1.1.1, Label B-L1, nexthop self, and transport-target:0:100. MPLS swap route is installed at ASBR13 for B-L1 with a nexthop pointing to ASBR13_to_PE11_gold tunnel.

Similarly ASBR14 also receives BGP CT route for 1.1.1.1:10:1.1.1.1 from PE11 and it resolves over the tunnel ASBR14_to_PE11_gold. The route is then readvertised by ASBR14 in BGP CT family to ASBRs ASBR21, ASBR22 according to export policy. This route is sent with same NLRI RD prefix 1.1.1.1:10:1.1.1.1, Label B-L2, nexthop self, and transport-target:0:100. MPLS swap route is installed at ASBR14 for B-L1 with a nexthop pointing to ASBR14_to_PE11_gold tunnel.

In the Bronze plane, BGP CT route with Bronze SLA to endpoint PE11 is originated by PE11 with a NLRI containing RD prefix 1.1.1.1:20:1.1.1.1, and appropriate label. The RD allows both Gold and Bronze advertisements traverse path selection pinchpoints without any path hiding at RRs or ASBRs. And route target extended community transport-target:0:200 lets the route resolve over Bronze tunnels in the network, similar to the process being described for Gold SLA path.

Moving back to the Gold plane, ASBR21 receives the Gold SLA BGP CT routes for NLRI RD prefix 1.1.1.1:10:1.1.1.1 over the single hop EBGp sessions from ASBR13, ASBR14, and can compute ECMP/FRR towards them. ASBR21 readvertises BGP CT route for 1.1.1.1:10:1.1.1.1 with nexthop self (loopback address 2.2.2.1) to RR27, advertising a new label B-L3. MPLS swap route is installed for label B-L3 at ASBR21 to swap to received label B-L1, B-L2 and forward to ASBR13, ASBR14 respectively. RR27 readvertises this BGP CT route to ABR23, ABR24 with label and nexthop unchanged.

Similarly, ASBR22 receives BGP CT route 1.1.1.1:10:1.1.1.1 over the single hop EBGp sessions from ASBR13, ASBR14, and readvertises with nexthop self (loopback address 2.2.2.2) to RR27, advertising a new label B-L4. MPLS swap route is installed for label B-L4 at ASBR22 to swap to received label B-L1, B-L2 and forward to ASBR13, ASBR14 respectively. RR27 readvertises this BGP CT route also to ABR23, ABR24 with label and nexthop unchanged.

Addpath is enabled for BGP CT family on the sessions between RR27 and ASBRs, ABRs such that routes for 1.1.1.1:10:1.1.1.1 with the nexthops ASBR21 and ASBR22 are reflected to ABR23, ABR24 without any path hiding. Thus giving ABR23 visibility of both available nexthops for Gold SLA.

ABR23 receives the route with nexthop 2.2.2.1, label B-L3 from RR27. The route target "transport-target:0:100" on this route acts as Mapping Community, and instructs ABR23 to strictly resolve the nexthop using transport class 100 routes only. ABR23 is unable to find a route for 2.2.2.1 with transport class 100. Thus it considers this route unusable and does not propagate it further. This prunes ASBR21 from Gold SLA tunneled path.

ABR23 also receives the route with nexthop 2.2.2.2, label B-L4 from RR27. The route target "transport-target:0:100" on this route acts as Mapping Community, and instructs ABR23 to strictly resolve the nexthop using transport class 100 routes only. ABR23 successfully resolves the nexthop to point to ABR23_to_ASBR22_gold tunnel. ABR23 readvertises this BGP CT route with nexthop self (loopback address 2.2.2.3) and a new label B-L5 to PE25. Swap route for B-L5 is installed by ABR23 to swap to label B-L4, and forward into ABR23_to_ASBR22_gold tunnel.

PE25 receives the BGP CT route for prefix 1.1.1.1:10:1.1.1.1 with label B-L5, nexthop 2.2.2.3 and transport-target:0:100 from RR26. And it similarly resolves the nexthop 2.2.2.3 over transport class 100, pushing labels associated with PE25_to_ABR23_gold tunnel.

In this manner, the Gold transport LSP "ASBR13_to_PE11_gold" in egress-domain is extended by BGP CT until the ingress-node PE25 in ingress domain, to create an end-to-end Gold SLA path. MPLS swap routes are installed at ASBR13, ASBR22 and ABR23, when propagating the PE11 BGP CT Gold transport class route 1.1.1.1:10:1.1.1.1 with nexthop self towards PE25.

The BGP CT LSP thus formed, originates in PE25, and terminates in ASBR13 (assuming PE11 advertised Implicit Null), traversing over the Gold underlay LSPs in each domain. ASBR13 uses UHP to stitch the BGP CT LSP into the "ASBR13_to_PE11_gold" LSP to traverse the last domain, thus satisfying Gold SLA end-to-end.

When PE25 receives service routes from RR26 with nexthop 1.1.1.1 and mapping community Color:0:100, it resolves over this BGP CT route 1.1.1.1:10:1.1.1.1. Thus pushing label B-L5, and pushing as top label the labels associated with PE25_to_ABR23_gold tunnel.

18.4. Data plane view

18.4.1. Steady state

This section describes how the data plane looks like in steady state.

CE41 transmits an IP packet with destination as 31.31.31.31. On receiving this packet PE25 performs a lookup in the IP FIB associated with the CE41 interface. This lookup yields the service route that pushes the VPN service label V-L1, BGP CT label B-L5, and labels for PE25_to_ABR23_gold tunnel. Thus PE25 encapsulates the IP packet in MPLS packet with label V-L1(innermost), B-L5, and top label as PE25_to_ABR23_gold tunnel. This MPLS packet is thus transmitted to ABR23 using Gold SLA.

ABR23 decapsulates the packet received on PE25_to_ABR23_gold tunnel as required, and finds the MPLS packet with label B-L5. It performs lookup for label B-L5 in the global MPLS FIB. This yields the route that swaps label B-L5 with label B-L4, and pushes top label provided by ABR23_to_ASBR22_gold tunnel. Thus ABR23 transmits the MPLS packet with label B-L4 to ASBR22, on a tunnel that satisfies Gold SLA.

ASBR22 similarly performs a lookup for label B-L4 in global MPLS FIB, finds the route that swaps label B-L4 with label B-L2, and forwards to ASBR13 over the directly connected MPLS enabled interface. This interface is a common resource not dedicated to any specific transport class, in this example.

ASBR13 receives the MPLS packet with label B-L2, and performs a lookup in MPLS FIB, finds the route that pops label B-L2, and pushes labels associated with ASBR13_to_PE11_gold tunnel. This transmits the MPLS packet with VPN label V-L1 to PE11 using a tunnel that preserves Gold SLA in AS 1.

PE11 receives the MPLS packet with V-L1, and performs VPN forwarding. Thus transmitting the original IP payload from CE41 to CE31. The payload has traversed path satisfying Gold SLA end-to-end.

18.4.2. Local repair of primary path

This section describes how the data plane at ASBR22 reacts when link between ASBR22 and ASBR13 experiences a failure, and an alternate path exists.

Assuming ASBR22_to_ASBR13 link goes down, such that traffic with Gold SLA going to PE11 needs repair. ASBR22 has an alternate BGP CT route for 1.1.1.1:10:1.1.1.1 from ASBR14. This has been preprogrammed in forwarding by ASBR22 as FRR backup nexthop for label B-L4. This

allows the Gold SLA traffic to be locally repaired at ASBR22 without the failure event propagated in the BGP CT network. In this case, ingress node PE25 will not know there was a failure, and traffic restoration will be independent of prefix scale (PIC).

18.4.3. Absorbing failure of primary path. Fallback to best-effort tunnels.

This section describes how the data plane reacts when gold path experiences a failure, but no alternate path exists.

Assuming tunnel ABR23_to_ASBR22_gold goes down, such that now end-to-end Gold path does not exist in the network. This makes the BGP CT route for RD prefix 1.1.1.1:10:1.1.1.1 unusable at ABR23. This makes ABR23 send a BGP withdrawal for 1.1.1.1:10:1.1.1.1 to PE25.

Withdrawal for 1.1.1.1:10:1.1.1.1 allows PE25 to react to the loss of gold path to 1.1.1.1. Assuming PE25 is provisioned to use best-effort transport class as the backup path, this withdrawal of BGP CT route allows PE25 to adjust the nexthop of the VPN Service-route to push the labels provided by the BGP LU route. That repairs the traffic to go via best effort path. PE25 can also be provisioned to use Bronze transport class as the backup path. The repair will happen in similar manner in that case as-well.

Traffic repair to absorb the failure happens at ingress node PE25, in a service prefix scale independent manner. This is called PIC (Prefix scale Independent Convergence). The repair time will be proportional to time taken for withdrawing the BGP CT route.

The above examples demonstrate the various levels of failsafe mechanisms available to protect traffic in a BGP CT network.

19. IANA Considerations

This document makes following requests of IANA.

19.1. New BGP SAFI

New BGP SAFI code for "Classful Transport". Value 76.

This will be used to create new AFI,SAFI pairs for IPv4, IPv6 Classful Transport families. viz:

- * "Inet, Classful Transport". AFI/SAFI = "1/76" for carrying IPv4 Classful Transport prefixes.

- * "Inet6, Classful Transport". AFI/SAFI = "2/76" for carrying IPv6 Classful Transport prefixes.

19.2. New Format for BGP Extended Community

Please assign a new Format (Type high = 0xa) of extended community EXT-COMM [RFC4360] called "Transport Class" from the following registries:

the "BGP Transitive Extended Community Types" registry, and

the "BGP Non-Transitive Extended Community Types" registry.

Please assign the same low-order six bits for both allocations.

This document uses this new Format with subtype 0x2 (route target), as a transitive extended community.

The Route Target thus formed is called "Transport Class" route target extended community.

Taking reference of RFC7153 [RFC7153] , following requests are made:

19.2.1. Existing registries to be modified

19.2.1.1. Registries for the "Type" Field

19.2.1.1.1. Transitive Types

This registry contains values of the high-order octet (the "Type" field) of a Transitive Extended Community.

Registry Name: BGP Transitive Extended Community Types

	TYPE VALUE	NAME
+	0x0a	Transitive Transport Class Extended
+		Community (Sub-Types are defined in the
+		"Transitive Transport Class Extended
+		Community Sub-Types" registry)

19.2.1.1.2. Non-Transitive Types

This registry contains values of the high-order octet (the "Type" field) of a Non-transitive Extended Community.

Registry Name: BGP Non-Transitive Extended Community Types

TYPE	VALUE	NAME
+	0x4a	Non-Transitive Transport Class Extended
+		Community (Sub-Types are defined in the
+		"Non-Transitive Transport Class Extended
+		Community Sub-Types" registry)

19.2.2. New registries to be created

19.2.2.1. Transitive "Transport Class" Extended Community Sub-Types Registry

This registry contains values of the second octet (the "Sub-Type" field) of an extended community when the value of the first octet (the "Type" field) is 0x07.

Registry Name: Transitive Transport Class Extended Community Sub-Types

RANGE	REGISTRATION PROCEDURE
0x00-0xBF	First Come First Served
0xC0-0xFF	IETF Review
SUB-TYPE VALUE	NAME
0x02	Route Target

19.2.2.2. Non-Transitive "Transport Class" Extended Community Sub-Types Registry

This registry contains values of the second octet (the "Sub-Type" field) of an extended community when the value of the first octet (the "Type" field) is 0x47.

Registry Name: Non-Transitive Transport Class Extended Community Sub-Types

RANGE	REGISTRATION PROCEDURE
0x00-0xBF	First Come First Served
0xC0-0xFF	IETF Review
SUB-TYPE VALUE	NAME
0x02	Route Target

19.3. MPLS OAM code points

The following two code points are sought for Target FEC Stack sub-TLVs:

- * IPv4 BGP Classful Transport
- * IPv6 BGP Classful Transport

20. Security Considerations

Mechanisms described in this document carry Transport routes in a new BGP address family. That minimizes possibility of these routes leaking outside the expected domain or mixing with service routes.

When redistributing between SAFI 4 and SAFI 76 Classful Transport routes, there is a possibility of SAFI 4 routes mixing with SAFI 1 service routes. To avoid such scenarios, it is RECOMMENDED that implementations support keeping SAFI 4 routes in a separate transport RIB, distinct from service RIB that contain SAFI 1 service routes.

21. Contributors

Rajesh M
Juniper Networks, Inc.
Electra, Exora Business Park~Marathahalli - Sarjapur Outer Ring Road,
Bangalore 560103
KA
India
Email: mrajesh@juniper.net

22. Acknowledgements

The authors thank Jeff Haas, John Scudder, Susan Hares, Reshma Das, Navaneetha Krishnan, Ravi M R, Chandrasekar Ramachandran, Shradha Hegde, Richard Roberts, Krzysztof Szarkowicz, John E Drake, Srihari Sangli, Vijay Kestur, Santosh Kolenchery, Robert Raszuk, Ahmed Darwish, Aravind Srinivas Srinivasa Prabhakar, Moshiko Nayman, Chris Trip for the valuable discussions and review comments.

The decision to not reuse SAFI 128 and create a new address-family to carry these transport-routes was based on suggestion made by Richard Roberts and Krzysztof Szarkowicz.

23. Normative References

[BGP-LU-EPE]

Gredler, Ed., "Egress Peer Engineering using BGP-LU", 6 July 2021, <<https://datatracker.ietf.org/doc/html/draft-gredler-idr-bgplu-epe-14>>.

[FLOWSPEC-REDIR-IP]

Simpson, Ed., "BGP Flow-Spec Redirect to IP Action", 2 February 2015, <<https://datatracker.ietf.org/doc/html/draft-ietf-idr-flowspec-redirect-ip-02>>.

[MPLS-NAMESPACES]

Vairavakkalai, Ed., "BGP signalled MPLS-namespaces", 11 June 2021, <<https://tools.ietf.org/html/draft-kaliraj-bess-bgp-sig-private-mpls-labels-01#section-6.1>>.

[MULTI-NH-ATTR]

Vairavakkalai, Ed., "BGP MultiNexthop Attribute", 28 December 2021, <<https://datatracker.ietf.org/doc/html/draft-kaliraj-idr-multinexthop-attribute-02#section-3.4.3>>.

[PCEP-RSVP-COLOR]

Rajagopalan, Ed., "Path Computation Element Protocol (PCEP) Extension for RSVP Color", 15 January 2021, <<https://datatracker.ietf.org/doc/html/draft-rajagopalan-pcep-rsvp-color-00>>.

[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>>.

[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>>.

[RFC4364] Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)", RFC 4364, DOI 10.17487/RFC4364, February 2006, <<https://www.rfc-editor.org/info/rfc4364>>.

- [RFC4456] Bates, T., Chen, E., and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)", RFC 4456, DOI 10.17487/RFC4456, April 2006, <<https://www.rfc-editor.org/info/rfc4456>>.
- [RFC4684] Marques, P., Bonica, R., Fang, L., Martini, L., Raszuk, R., Patel, K., and J. Guichard, "Constrained Route Distribution for Border Gateway Protocol/MultiProtocol Label Switching (BGP/MPLS) Internet Protocol (IP) Virtual Private Networks (VPNs)", RFC 4684, DOI 10.17487/RFC4684, November 2006, <<https://www.rfc-editor.org/info/rfc4684>>.
- [RFC4760] Bates, T., Chandra, R., Katz, D., and Y. Rekhter, "Multiprotocol Extensions for BGP-4", RFC 4760, DOI 10.17487/RFC4760, January 2007, <<https://www.rfc-editor.org/info/rfc4760>>.
- [RFC7153] Rosen, E. and Y. Rekhter, "IANA Registries for BGP Extended Communities", RFC 7153, DOI 10.17487/RFC7153, March 2014, <<https://www.rfc-editor.org/info/rfc7153>>.
- [RFC7911] Walton, D., Retana, A., Chen, E., and J. Scudder, "Advertisement of Multiple Paths in BGP", RFC 7911, DOI 10.17487/RFC7911, July 2016, <<https://www.rfc-editor.org/info/rfc7911>>.
- [RFC8029] Kompella, K., Swallow, G., Pignataro, C., Ed., Kumar, N., Aldrin, S., and M. Chen, "Detecting Multiprotocol Label Switched (MPLS) Data-Plane Failures", RFC 8029, DOI 10.17487/RFC8029, March 2017, <<https://www.rfc-editor.org/info/rfc8029>>.
- [RFC8212] Mauch, J., Snijders, J., and G. Hankins, "Default External BGP (EBGP) Route Propagation Behavior without Policies", RFC 8212, DOI 10.17487/RFC8212, July 2017, <<https://www.rfc-editor.org/info/rfc8212>>.
- [RFC8277] Rosen, E., "Using BGP to Bind MPLS Labels to Address Prefixes", RFC 8277, DOI 10.17487/RFC8277, October 2017, <<https://www.rfc-editor.org/info/rfc8277>>.
- [RFC8664] Sivabalan, S., Filsfils, C., Tantsura, J., Henderickx, W., and J. Hardwick, "Path Computation Element Communication Protocol (PCEP) Extensions for Segment Routing", RFC 8664, DOI 10.17487/RFC8664, December 2019, <<https://www.rfc-editor.org/info/rfc8664>>.

- [RFC8669] Previdi, S., Filsfils, C., Lindem, A., Ed., Sreekantiah, A., and H. Gredler, "Segment Routing Prefix Segment Identifier Extensions for BGP", RFC 8669, DOI 10.17487/RFC8669, December 2019, <<https://www.rfc-editor.org/info/rfc8669>>.
- [RFC8986] Filsfils, C., Ed., Camarillo, P., Ed., Leddy, J., Voyer, D., Matsushima, S., and Z. Li, "Segment Routing over IPv6 (SRv6) Network Programming", RFC 8986, DOI 10.17487/RFC8986, February 2021, <<https://www.rfc-editor.org/info/rfc8986>>.
- [RFC9012] Patel, K., Van de Velde, G., Sangli, S., and J. Scudder, "The BGP Tunnel Encapsulation Attribute", RFC 9012, DOI 10.17487/RFC9012, April 2021, <<https://www.rfc-editor.org/info/rfc9012>>.
- [RTC-Ext] Zhang, Z., Ed., "Route Target Constrain Extension", 12 July 2020, <<https://tools.ietf.org/html/draft-zzhang-idr-bgp-rt-constrains-extension-00#section-2>>.
- [Seamless-SR] Hegde, Ed., "Seamless Segment Routing", 17 November 2020, <<https://datatracker.ietf.org/doc/html/draft-hegde-spring-mpls-seamless-sr-03>>.
- [SRTE] Previdi, S., Ed., "Advertising Segment Routing Policies in BGP", 18 November 2019, <<https://tools.ietf.org/html/draft-ietf-idr-segment-routing-te-policy-08>>.
- [SRV6-INTER-DOMAIN] K A, Ed., "SRv6 inter-domain mapping SIDs", 10 January 2021, <<https://datatracker.ietf.org/doc/html/draft-salih-spring-srv6-inter-domain-sids-00>>.
- [SRV6-MPLS-AGRWL] Agrawal, Ed., "SRv6 and MPLS interworking", 22 February 2021, <<https://datatracker.ietf.org/doc/draft-agrawal-spring-srv6-mpls-interworking/05/>>.
- [SRV6-SERVICES] Dawra, Ed., "SRv6 BGP based Overlay Services", 11 April 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-bess-srv6-services-07>>.

Authors' Addresses

Kaliraj Vairavakkalai (editor)
Juniper Networks, Inc.
1133 Innovation Way,
Sunnyvale, CA 94089
United States of America
Email: kaliraj@juniper.net

Natrajan Venkataraman
Juniper Networks, Inc.
1133 Innovation Way,
Sunnyvale, CA 94089
United States of America
Email: natv@juniper.net

Balaji Rajagopalan
Juniper Networks, Inc.
Electra, Exora Business Park~Marathahalli - Sarjapur Outer Ring Road,
Bangalore 560103
KA
India
Email: balajir@juniper.net

Gyan Mishra
Verizon Communications Inc.
13101 Columbia Pike
Silver Spring, MD 20904
United States of America
Email: gyan.s.mishra@verizon.com

Mazen Khaddam
Cox Communications Inc.
Atlanta, GA
United States of America
Email: mazen.khaddam@cox.com

Xiaohu Xu
Capitalonline.
Beijing
China
Email: xiaohu.xu@capitalonline.net

Rafal Jan Szarecki
Google.
1160 N Mathilda Ave, Bldg 5,
Sunnyvale,, CA 94089
United States of America
Email: szarecki@google.com

Deepak J Gowda
Extreme Networks
55 Commerce Valley Drive West, Suite 300,
Thornhill, Toronto, Ontario L3T 7V9
Canada
Email: dgowda@extremenetworks.com

Chaitanya Yadlapalli
AT&T
200 S Laurel Ave,
Middletown,, NJ 07748
United States of America
Email: cy098d@att.com

Israel Means
AT&T
2212 Avenida Mara,
Chula Vista,, California 91914
United States of America
Email: israel.means@att.com