### BGP Extensions for Routing Policy Distribution (RPD)
### draft-ietf-idr-rpd-00

Abstract

   It is hard to adjust traffic and optimize traffic paths on a
   traditional IP network from time to time through manual
   configurations.  It is desirable to have an automatic mechanism for
   setting up routing policies, which adjust traffic and optimize
   traffic paths automatically.  This document describes BGP Extensions
   for Routing Policy Distribution (BGP RPD) to support this.

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 will expire on May 4, 2020.

Copyright Notice

Table of Contents

## 1.  Introduction

   It is difficult to optimize traffic paths on a traditional IP network
   because of:

   o  Heavy configuration and error prone.  Traffic can only be adjusted
      device by device.  All routers that the traffic traverses need to
      be configured.  The configuration workload is heavy.  The

operation is not only time consuming but also prone to
misconfiguration for Service Providers.

o  Complex.  The routing policies used to control network routes are
   complex, posing difficulties to subsequent maintenance, high
   maintenance skills are required.

It is desirable to have an automatic mechanism for setting up routing
policies, which can simplify the routing policies configuration.
This document describes extensions to BGP for Routing Policy
Distribution to resolve these issues.

## 2.  Terminology

The following terminology is used in this document.

o  ACL:Access Control List

o  BGP: Border Gateway Protocol

o  FS: Flow Specification

o  PBR:Policy-Based Routing

o  RPD: Routing Policy Distribution

o  VPN: Virtual Private Network

## 3.  Problem Statements

It is obvious that providers have the requirements to adjust their
business traffic from time to time because:

o  Business development or network failure introduces link congestion
   and overload.

o  Network transmission quality is decreased as the result of delay,
   loss and they need to adjust traffic to other paths.

o  To control OPEX and CPEX, prefer the transit provider with lower
   price.

### 3.1.  Inbound Traffic Control

In the scenario below, for the reasons above, the provider of AS100
saying P may wish the inbound traffic from AS200 enters AS100 through
link L3 instead of the others.  Since P doesn't have any

administration over AS200, so there is no way for P to modify the
route selection criteria directly.

```
                  Traffic from PE1 to Prefix1
             ----------------------------------->

+----------------+              +------------------------+
|     +---------+ |        L1  | +----+      +----------+|
|     |Speaker1 | +------------+ |IGW1|      |policy    ||
|     +---------+ |**     L2**| +----+      |controller||
|                 |  **     ** |              +----------+|
| +---+           |    ****    |                         |
| |PE1|           |    ****    |                         |
| +---+           |  **     ** |                         |
|     +---------+ |**     L3**| +----+                    |
|     |Speaker2 | +------------+ |IGW2|      AS100         |
|     +---------+ |      L4  | +----+                     |
|                 |            |                         |
|     AS200       |            |                         |
|                 |            |  ...                     |
|                 |            |                         |
|     +---------+ |            | +----+      +-------+    |
|     |Speakern | |            | |IGWn|      |Prefix1|    |
|     +---------+ |            | +----+      +-------+    |
+----------------+              +------------------------+

          Prefix1 advertised from AS100 to AS200
       <---------------------------------------

             Inbound Traffic Control case
```

## 3.2.  Outbound Traffic Control

In the scenario below, the provider of AS100 saying P prefers link L3
for the traffic to the destination Prefix2 among multiple exits and
links.  This preference can be dynamic and changed frequently because
of the reasons above.  So the provider P expects an efficient and
convenient solution.

```
                  Traffic from PE2 to Prefix2
           ------------------------------------>
   +-------------------------+          +----------------+
   |+----------+     +----+ |L1          | +---------+      |
   ||policy    |     |IGW1| +------------+ |Speaker1 |      |
   ||controller|     +----+ |**        **| +---------+      |
   |+----------+            |L2**     ** |        +-------+|
   |                        |    ****    |        |Prefix2||
   |                        |    ****    |        +-------+|
   |                        |L3**     ** |                 |
   |       AS100            +----+ |**        **| +---------+      |
   |                        |IGW2| +------------+ |Speaker2 |      |
   |                        +----+ |L4          | +---------+      |
   |                               |            |                 |
   |+---+                          |            |     AS200        |
   ||PE2|            ...    |            |                 |
   |+---+                          |            |                 |
   |                        +----+ |            | +---------+      |
   |                        |IGWn| |            | |Speakern |      |
   |                        +----+ |            | +---------+      |
   +------------------------+          +----------------+

            Prefix2 advertised from AS200 to AS100
           <---------------------------------------

                  Outbound Traffic Control case
```

## 4.  Protocol Extensions

A solution is proposed to use a new AFI and SAFI with the BGP Wide
Community for encoding a routing policy.

### 4.1.  Using a New AFI and SAFI

A new AFI and SAFI are defined: the Routing Policy AFI whose
codepoint TBD1 is to be assigned by IANA, and SAFI whose codepoint
TBD2 is to be assigned by IANA.

The AFI and SAFI pair uses a new NLRI, which is defined as follows:

```
 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
+-+-+-+-+-+-+-+-+
|  NLRI Length  |
+-+-+-+-+-+-+-+-+
|  Policy Type  |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Distinguisher (4 octets)                 |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Peer IP (4/16 octets)                    ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

Where:

  NLRI Length:  1 octet represents the length of NLRI.

  Policy Type:  1 octet indicates the type of a policy.  1 is for
   export policy. 2 is for import policy.

  Distinguisher:  4 octet value uniquely identifies the policy in the
   peer.

  Peer IP:  4/16 octet value indicates an IPv4/IPv6 peer.

The NLRI containing the Routing Policy is carried in a BGP UPDATE
message, which MUST contain the BGP mandatory attributes and MAY also
contain some BGP optional attributes.

When receiving a BGP UPDATE message, a BGP speaker processes it only
if the peer IP address in the NLRI is the IP address of the BGP
speaker or 0.

The content of the Routing Policy is encoded in a BGP Wide Community.

## 4.2.  BGP Wide Community

The BGP wide community is defined in
[I-D.ietf-idr-wide-bgp-communities].  It can be used to facilitate
the delivery of new network services, and be extended easily for
distributing different kinds of routing policies.

### 4.2.1.  New Wide Community Atoms

A wide community Atom is a TLV (or sub-TLV), which may be included in
a BGP wide community container (or BGP wide community for short)
containing some BGP Wide Community TLVs.  Three BGP Wide Community
TLVs are defined in [I-D.ietf-idr-wide-bgp-communities], which are
BGP Wide Community Target(s) TLV, Exclude Target(s) TLV, and

Parameter(s) TLV.  Each of these TLVs comprises a series of Atoms,
each of which is a TLV (or sub-TLV).  A new wide community Atom is
defined for BGP Wide Community Target(s) TLV and a few new Atoms are
defined for BGP Wide Community Parameter(s) TLV.  For your reference,
the format of the TLV is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-++-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Type      |              Length           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                      Value (variable)                       ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                 Format of Wide Community Atom TLV

A RouteAttr Atom TLV (or RouteAttr TLV/sub-TLV for short) is defined
and may be included in a Target TLV.  It has the following format.

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Type (TBD1)   |       Length (variable)       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         sub-TLVs                            ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                 Format of RouteAttr Atom TLV

The Type for RouteAttr is TBD1 (suggested value 48) to be assigned by
IANA.  In RouteAttr TLV, three sub-TLVs are defined: IP Prefix, AS-
Path and Community sub-TLV.

An IP prefix sub-TLV gives matching criteria on IPv4 prefixes.  Its
format is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Type (TBD2)  |          Length (N x 8)       |M-Type | Flags |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         IPv4 Address                         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Mask     |     GeMask    |     LeMask    |M-Type | Flags |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
~       . . .
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         IPv4 Address                         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Mask     |     GeMask    |     LeMask    |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                   Format of IPv4 Prefix sub-TLV

   Type:  TBD2 (suggested value 1) for IPv4 Prefix is to be assigned by
      IANA.

   Length:  N x 8, where N is the number of tuples <M-Type, Flags, IPv4
      Address, Mask, GeMask, LeMask>.

   M-Type:  4 bits for match types, four of which are defined:

      M-Type = 0:  Exact match.

      M-Type = 1:  Match prefix greater and equal to the given masks.

      M-Type = 2:  Match prefix less and equal to the given masks.

      M-Type = 3:  Match prefix within the range of the given masks.

   Flags:  4 bits.  No flags are currently defined.

   IPv4 Address:  4 octets for an IPv4 address.

   Mask:  1 octet for the mask length.

   GeMask:  1 octet for match range, must be less than Mask or be 0.

   LeMask:  1 octet for match range, must be greater than Mask or be 0.

   For example, tuple <M-Type=0, Flags=0, IPv4 Address = 1.1.0.0, Mask =
   22, GeMask = 0, LeMask = 0> represents an exact IP prefix match for
   1.1.0.0/22.

   <M-Type=1, Flags=0, IPv4 Address = 16.1.0.0, Mask = 24, GeMask = 24,
   LeMask = 0> represents match IP prefix 1.1.0.0/24 greater-equal 24.

   <M-Type=2, Flags=0, IPv4 Address = 17.1.0.0, Mask = 24, GeMask = 0,
   LeMask = 26> represents match IP prefix 17.1.0.0/24 less-equal 26.

   <M-Type=3, Flags=0, IPv4 Address = 18.1.0.0, Mask = 24, GeMask = 24,
   LeMask = 32> represents match IP prefix 18.1.0.0/24 greater-equal to
   24 and less-equal 32.

   Similarly, an IPv6 Prefix sub-TLV represents match criteria on IPv6
   prefixes.  Its format is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Type(TBD3)  |          Length (N x 20)      |M-Type | Flags |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     IPv6 Address (16 octets)                ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Mask      |     GeMask    |     LeMask    |M-Type | Flags |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
~      . . .
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     IPv6 Address (16 octets                 ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Mask      |     GeMask    |     LeMask    |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                   Format of IPv6 Prefix sub-TLV

   An AS-Path sub-TLV represents a match criteria in a regular
   expression string.  Its format is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Type (TBD4)  |      Length (Variable)       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     AS-Path Regex String                    |
:                                                             :
|                                                             ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                    Format of AS Path sub-TLV

   Type:  TBD4 (suggested value 2) for AS-Path is to be assigned by
      IANA.

Length:  Variable, maximum is 1024.

AS-Path Regex String:  AS-Path regular expression string.

A community sub-TLV represents a list of communities to be matched
all.  Its format is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|  Type (TBD5)  |        Length (N x 4 + 1)     |    Flags      |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Community 1 Value                         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
~                         . . .                                 ~
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Community N Value                         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                    Format of Community sub-TLV

Type:  TBD5 (suggested value 3) for Community is to be assigned by
   IANA.

Length:  N x 4 + 1, where N is the number of communities.

Flags:  1 octet.  No flags are currently defined.

In Parameter(s) TLV, two action sub-TLVs are defined: MED change sub-
TLV and AS-Path change sub-TLV.  When the community in the container
is MATCH AND SET ATTR, the Parameter(s) TLV includes some of these
sub-TLVs.  When the community is MATCH AND NOT ADVERTISE, the
Parameter(s) TLV's value is empty.

A MED change sub-TLV indicates an action to change the MED.  Its
format is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|  Type (TBD6)  |         Length (5)            |      OP       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                        Value                                  |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

                    Format of MED Change sub-TLV

Type:  TBD6 (suggested value 1) for MED Change is to be assigned by
   IANA.

Length:  5.

OP:  1 octet.  Three are defined:
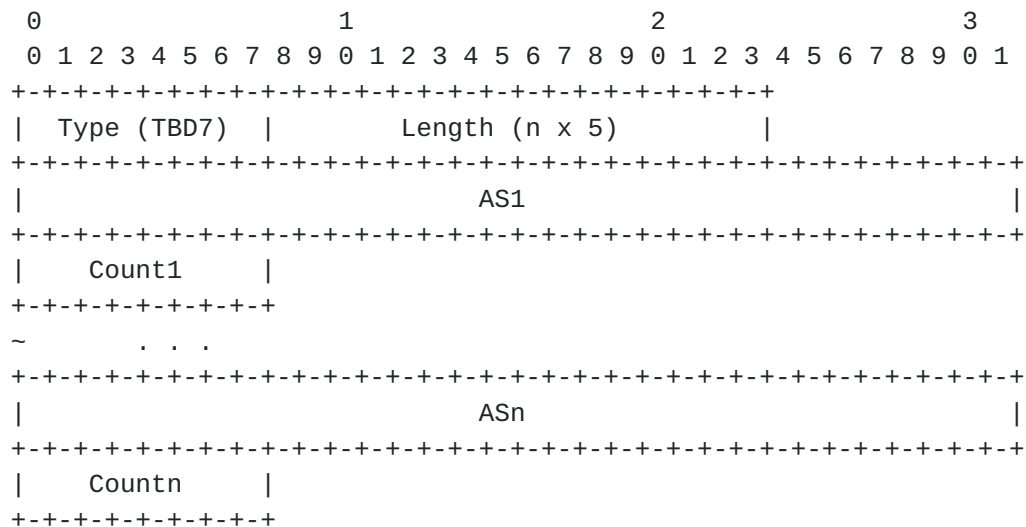
   OP = 0:  assign the Value to the existing MED.

   OP = 1:  add the Value to the existing MED.  If the sum is greater
      than the maximum value for MED, assign the maximum value to
      MED.

   OP = 2:  subtract the Value from the existing MED.  If the
      existing MED minus the Value is less than 0, assign 0 to MED.

Value:  4 octets.

An AS-Path change sub-TLV indicates an action to change the AS-Path.
Its format is illustrated below:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Type (TBD7) |        Length (n x 5)          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                            AS1                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Count1     |
+-+-+-+-+-+-+-+-+
~       . . .
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                            ASn                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Countn     |
+-+-+-+-+-+-+-+-+
```

                  Format of AS-Path Change sub-TLV

Type:  TBD7 (suggested value 2) for AS-Path Change is to be assigned
   by IANA.

Length:  n x 5.

ASi:  4 octet.  An AS number.

Counti:  1 octet.  ASi repeats Counti times.

The sequence of AS numbers are added to the existing AS Path.

## 4.3.  Capability Negotiation

It is necessary to negotiate the capability to support BGP Extensions
for Routing Policy Distribution (RPD).  The BGP RPD Capability is a
new BGP capability [RFC5492].  The Capability Code for this
capability is to be specified by the IANA.  The Capability Length
field of this capability is variable.  The Capability Value field
consists of one or more of the following tuples:

```
+---------------------------------------------------+
|  Address Family Identifier (2 octets)             |
+---------------------------------------------------+
|  Subsequent Address Family Identifier (1 octet)   |
+---------------------------------------------------+
|  Send/Receive (1 octet)                           |
+---------------------------------------------------+
```

BGP RPD Capability

The meaning and use of the fields are as follows:

Address Family Identifier (AFI): This field is the same as the one
used in [RFC4760].

Subsequent Address Family Identifier (SAFI): This field is the same
as the one used in [RFC4760].

Send/Receive: This field indicates whether the sender is (a) willing
to receive Routing Policies from its peer (value 1), (b) would like
to send Routing Policies to its peer (value 2), or (c) both (value 3)
for the <AFI, SAFI>.

## 5.  Consideration

## 5.1.  Route-Policy

Routing policies are used to filter routes and control how routes are
received and advertised.  If route attributes, such as reachability,
are changed, the path along which network traffic passes changes
accordingly.

When advertising, receiving, and importing routes, the router
implements certain policies based on actual networking requirements
to filter routes and change the attributes of the routes.  Routing
policies serve the following purposes:

o  Control route advertising: Only routes that match the rules
   specified in a policy are advertised.

o  Control route receiving: Only the required and valid routes are
   received.  This reduces the size of the routing table and improves
   network security.

o  Filter and control imported routes: A routing protocol may import
   routes discovered by other routing protocols.  Only routes that
   satisfy certain conditions are imported to meet the requirements
   of the protocol.

o  Modify attributes of specified routes Attributes of the routes:
   that are filtered by a routing policy are modified to meet the
   requirements of the local device.

o  Configure fast reroute (FRR): If a backup next hop and a backup
   outbound interface are configured for the routes that match a
   routing policy, IP FRR, VPN FRR, and IP+VPN FRR can be
   implemented.

   Routing policies are implemented using the following procedures:

1.  Define rules: Define features of routes to which routing policies
    are applied.  Users define a set of matching rules based on
    different attributes of routes, such as the destination address
    and the address of the router that advertises the routes.

2.  Implement the rules: Apply the matching rules to routing policies
    for advertising, receiving, and importing routes.

## 6.  Contributors

   The following people have substantially contributed to the definition
   of the BGP-FS RPD and to the editing of this document:

   Peng Zhou
   Huawei
   Email: Jewpon.zhou@huawei.com

## 7.  Security Considerations

   Protocol extensions defined in this document do not affect the BGP
   security other than those as discussed in the Security Considerations
   section of [RFC5575].

## 8.  Acknowledgements

   The authors would like to thank Acee Lindem, Jeff Haas, Jie Dong,
   Lucy Yong, Qiandeng Liang, Zhenqiang Li for their comments to this
   work.

## 9.  IANA Considerations

   This document requests assigning a new AFI in the registry "Address
   Family Numbers" as follows:

```
   +----------------------+------------------------+-------------+
   | Code Point           | Description            | Reference   |
   +----------------------+------------------------+-------------+
   | TBD (36879 suggested) |  Routing Policy AFI    |This document|
   +----------------------+------------------------+-------------+
```

   This document requests assigning a new SAFI in the registry
   "Subsequent Address Family Identifiers (SAFI) Parameters" as follows:

```
   +----------------------+------------------------+-------------+
   | Code Point           | Description            | Reference   |
   +----------------------+------------------------+-------------+
   | TBD(179 suggested)   |  Routing Policy SAFI   |This document|
   +----------------------+------------------------+-------------+
```

   This document defines a new registry called "Routing Policy NLRI".
   The allocation policy of this registry is "First Come First Served
   (FCFS)" according to [RFC8126].

   Following code points are defined:

```
   +-------------+-----------------------------------+-------------+
   | Code Point  | Description                       | Reference   |
   +-------------+-----------------------------------+-------------+
   |     1       | Export Policy                     |This document|
   +-------------+-----------------------------------+-------------+
   |     2       | Import Policy                     |This document|
   +-------------+-----------------------------------+-------------+
```

   This document requests assigning a code-point from the registry "BGP
   Community Container Atom Types" as follows:

```
   +--------------------+----------------------------+-------------+
   | TLV Code Point     | Description                | Reference   |
   +--------------------+----------------------------+-------------+
   | TBD1 (48 suggested) | RouteAttr Atom            |This document|
   +--------------------+----------------------------+-------------+
```

This document defines a new registry called "Route Attributes Sub-TLV" under RouteAttr Atom TLV.  The allocation policy of this registry is "First Come First Served (FCFS)" according to [RFC8126].

Following Sub-TLV code points are defined:

```
+-------------+--------------------------------+-------------+
| Code Point  | Description                    | Reference   |
+-------------+--------------------------------+-------------+
|      0      |  Reserved                      |             |
+-------------+--------------------------------+-------------+
|      1      |  IP Prefix Sub-TLV             |This document|
+-------------+--------------------------------+-------------+
|      2      |  AS-Path Sub-TLV               |This document|
+-------------+--------------------------------+-------------+
|      3      |  Community Sub-TLV             |This document|
+-------------+--------------------------------+-------------+
|   4 - 255   |  To be assigned in FCFS        |             |
+-------------+--------------------------------+-------------+
```

This document defines a new registry called "Attribute Change Sub-TLV" under Parameter(s) TLV.  The allocation policy of this registry is "First Come First Served (FCFS)" according to [RFC8126].

Following Sub-TLV code points are defined:

```
+-------------+--------------------------------+-------------+
| Code Point  | Description                    | Reference   |
+-------------+--------------------------------+-------------+
|      0      |  Reserved                      |             |
+-------------+--------------------------------+-------------+
|      1      |  MED Change Sub-TLV            |This document|
+-------------+--------------------------------+-------------+
|      2      |  AS-Path Change Sub-TLV        |This document|
+-------------+--------------------------------+-------------+
|   3 - 255   |  To be assigned in FCFS        |             |
+-------------+--------------------------------+-------------+
```

## 10.  References

## 10.1.  Normative References

[I-D.ietf-idr-wide-bgp-communities]
          Raszuk, R., Haas, J., Lange, A., Decraene, B., Amante, S.,
          and P. Jakma, "BGP Community Container Attribute", draft-ietf-idr-wide-bgp-communities-05 (work in progress), July
          2018.

   [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>.

   [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>.

   [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>.

   [RFC5492]  Scudder, J. and R. Chandra, "Capabilities Advertisement
              with BGP-4", RFC 5492, DOI 10.17487/RFC5492, February
              2009, <https://www.rfc-editor.org/info/rfc5492>.

   [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>.

   [RFC8126]  Cotton, M., Leiba, B., and T. Narten, "Guidelines for
              Writing an IANA Considerations Section in RFCs", BCP 26,
              RFC 8126, DOI 10.17487/RFC8126, June 2017,
              <https://www.rfc-editor.org/info/rfc8126>.

## 10.2.  Informative References

   [I-D.ietf-idr-registered-wide-bgp-communities]
              Raszuk, R. and J. Haas, "Registered Wide BGP Community
              Values", draft-ietf-idr-registered-wide-bgp-communities-02
              (work in progress), May 2016.

Authors' Addresses

   Zhenbin Li
   Huawei
   Huawei Bld., No.156 Beiqing Rd.
   Beijing  100095
   China

   Email: lizhenbin@huawei.com


   Liang Ou
   China Telcom Co., Ltd.
   109 West Zhongshan Ave,Tianhe District
   Guangzhou  510630
   China

   Email: oul@gsta.com


   Yujia Luo
   China Telcom Co., Ltd.
   109 West Zhongshan Ave,Tianhe District
   Guangzhou  510630
   China

   Email: luoyuj@gsta.com


   Sujian Lu
   Tencent
   Tengyun Building,Tower A ,No. 397 Tianlin Road
   Shanghai, Xuhui District  200233
   China

   Email: jasonlu@tencent.com


   Huaimo Chen
   Futurewei
   Boston, MA
   USA

   Email: Huaimo.chen@futurewei.com

    Shunwan Zhuang
    Huawei
    Huawei Bld., No.156 Beiqing Rd.
    Beijing  100095
    China

    Email: zhuangshunwan@huawei.com


    Haibo Wang
    Huawei
    Huawei Bld., No.156 Beiqing Rd.
    Beijing  100095
    China

    Email: rainsword.wang@huawei.com