

Network Working Group  
Internet Draft  
Intended status: Standard  
Expires: August 21, 2021

L. Dunbar  
Futurewei  
S. Hares  
Hickory Hill Consulting  
R. Raszuk  
Bloomberg LP  
K. Majumdar  
CommScope  
February 21, 2021

**BGP UPDATE for SDWAN Edge Discovery**  
**draft-dunbar-idr-sdwan-edge-discovery-02**

Abstract

The document describes encoding of BGP UPDATE messages for the SDWAN edge node discovery.

In the context of this document, BGP Route Reflectors (RR) is the component of the SDWAN Controller that receives the BGP UPDATE from SDWAN edges and in turns propagates the information to the intended peers that are authorized to communicate via the SDWAN overlay network.

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), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>

This Internet-Draft will expire on Dec 21, 2021.

## Copyright Notice

Copyright (c) 2021 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the [Trust Legal Provisions](#) and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

<a href="#">1. Introduction.....</a>	<a href="#">3</a>
<a href="#">2. Conventions used in this document.....</a>	<a href="#">3</a>
<a href="#">3. Framework of SDWAN Edge Discovery.....</a>	<a href="#">4</a>
<a href="#">3.1. The Objectives of SDWAN Edge Discovery.....</a>	<a href="#">4</a>
<a href="#">3.2. Comparing with Pure IPsec VPN.....</a>	<a href="#">5</a>
<a href="#">3.3. Client UPDATE and Hybrid Underlay Tunnel UPDATE.....</a>	<a href="#">6</a>
<a href="#">3.4. Edge Node Discovery.....</a>	<a href="#">8</a>
<a href="#">4. BGP UPDATE to Support SDWAN Segmentation.....</a>	<a href="#">9</a>
4.1. SDWAN Segmentation, SDWAN Virtual Topology and Client VPN.9	
<a href="#">4.2. Constrained Propagation of Edge Capability.....</a>	<a href="#">10</a>
<a href="#">4.3. SDWAN VPN ID in BGP Update.....</a>	<a href="#">11</a>
<a href="#">4.4. SDWAN VPN ID in Data Plane.....</a>	<a href="#">13</a>
<a href="#">5. Hybrid Underlay Tunnel UPDATE.....</a>	<a href="#">13</a>
<a href="#">5.1. NLRI for Hybrid Underlay Tunnel Update.....</a>	<a href="#">13</a>
<a href="#">5.2. SDWAN-Hybrid Tunnel Encoding.....</a>	<a href="#">15</a>
<a href="#">5.3. IPsec-SA-ID Sub-TLV.....</a>	<a href="#">15</a>
<a href="#">5.3.1. Encoding example #1 of using IPsec-SA-ID Sub-TLV....</a>	<a href="#">17</a>
<a href="#">5.3.2. Encoding Example #2 of using IPsec-SA-ID.....</a>	<a href="#">18</a>
<a href="#">5.4. Extended Port Sub-TLV.....</a>	<a href="#">19</a>
<a href="#">5.5. ISP of the Underlay network Sub-TLV.....</a>	<a href="#">21</a>

<a href="#">6.</a>	<a href="#">IPsec SA Property Sub-TLVs.....</a>	<a href="#">23</a>
<a href="#">6.1.</a>	<a href="#">IPsec SA Nonce Sub-TLV.....</a>	<a href="#">23</a>
<a href="#">6.2.</a>	<a href="#">IPsec Public Key Sub-TLV.....</a>	<a href="#">24</a>
<a href="#">6.3.</a>	<a href="#">IPsec SA Proposal Sub-TLV.....</a>	<a href="#">24</a>
<a href="#">6.4.</a>	<a href="#">Simplified IPsec Security Association sub-TLV.....</a>	<a href="#">24</a>
<a href="#">6.5.</a>	<a href="#">IPsec SA Encoding Examples.....</a>	<a href="#">26</a>
<a href="#">7.</a>	<a href="#">Error &amp; Mismatch Handling.....</a>	<a href="#">27</a>
<a href="#">8.</a>	<a href="#">Manageability Considerations.....</a>	<a href="#">28</a>
<a href="#">9.</a>	<a href="#">Security Considerations.....</a>	<a href="#">28</a>
<a href="#">10.</a>	<a href="#">IANA Considerations.....</a>	<a href="#">28</a>
<a href="#">11.</a>	<a href="#">References.....</a>	<a href="#">29</a>
<a href="#">11.1.</a>	<a href="#">Normative References.....</a>	<a href="#">29</a>
<a href="#">11.2.</a>	<a href="#">Informative References.....</a>	<a href="#">29</a>
<a href="#">12.</a>	<a href="#">Acknowledgments.....</a>	<a href="#">30</a>

## [1.](#) Introduction

[SDWAN-BGP-USAGE] illustrates how BGP is used as control plane for a SDWAN network. SDWAN network refers to a policy-driven network over multiple different underlay networks to get better WAN bandwidth management, visibility and control.

The document describes a BGP UPDATE for SDWAN edge nodes to announce its properties to its RR which then propagates to the authorized peers.

## [2.](#) Conventions used in this document

Cloud DC: Off-Premise Data Centers that usually host applications and workload owned by different organizations or tenants.

Controller: Used interchangeably with SDWAN controller to manage SDWAN overlay path creation/deletion and monitor the path conditions between sites.

CPE-Based VPN: Virtual Private Secure network formed among CPEs. This is to differentiate from most commonly used PE-based VPNs a la [RFC 4364](#).

MP-NLRI: The MP\_REACH\_NLRI Path Attribute defined in [RFC4760](#).

SDWAN End-point: can be the SDWAN edge node address, a WAN port address (logical or physical) of a SDWAN edge node, or a client port address.

OnPrem: On Premises data centers and branch offices

SDWAN: Software Defined Wide Area Network. In this document, "SDWAN" refers to policy-driven transporting IP packets over multiple different underlay networks to get better WAN bandwidth management, visibility and control.

SDWAN Segmentation: Segmentation is the process of dividing the network into logical sub-networks.

SDWAN VPN: referring to Client's VPN, which is like the VRF on the PEs of a MPLS VPN. One SDWAN client VPN can be mapped one or multiple SD-WAN virtual topologies. How Client VPN is mapped to a SDWAN virtual topology is governed by policies.

SDWAN Virtual Topology: Since SDWAN can connect any nodes, whereas MPLS VPN connects a fixed number of PEs, one SDWAN Virtual Topology refers to a set of edge nodes and the tunnels (including both IPsec tunnels and/or MPLS tunnels) interconnecting those edge nodes.

### **3. Framework of SDWAN Edge Discovery**

#### **3.1. The Objectives of SDWAN Edge Discovery**

The objectives of SDWAN edge discovery is for a SDWAN edge node to discover its authorized peers and their associated properties for its attached clients traffic to communicate. The attributes to be propagated includes the SDWAN (client) VPNs supported, the attached routes under specific SDWAN VPNs, and the properties of the underlay networks over which the client routes can be carried.

Some SDWAN peers are connected by both trusted VPNs and untrusted public networks. Some SDWAN peers are connected only by untrusted public networks. For the portion over untrusted networks, IPsec Security Associations (IPsec SA) have to be established and

maintained. If an edge node has network ports behind the NAT, the NAT properties needs to be discovered by authorized SDWAN peers.

Just like any VPN networks, the attached client's routes belonging to specific SDWAN VPNs can only be exchanged to the SDWAN peer nodes that are authorized to communicate.

### **3.2. Comparing with Pure IPsec VPN**

Pure IPsec VPN has IPsec tunnels connecting all edge nodes via public internet, therefore requires stringent authentication and authorization (i.e. IKE Phase 1) before other properties of IPsec SA can be exchanged. The IPsec Security Association (SA) between two untrusted nodes typically requires the following configurations and message exchanges:

- IPsec IKE to authenticate with each other
- Establish IPsec SA
  - o Local key configuration
  - o Remote Peer address (192.10.0.10<->172.0.01)
  - o IKEv2 Proposal directly sent to peer
    - o Encryption method, Integrity sha512
  - o Transform set
- Attached client prefixes discovery
  - o By running routing protocol within each IPsec SA
  - o If multiple IPsec SAs between two peer nodes are established to achieve load sharing, each IPsec tunnel needs to run its own routing protocol to exchange client routes attached to the edges.
- Access List or Traffic Selector)
  - o Permit Local-IP1, Remote-IP2

In a BGP controlled SDWAN network, e.g. a MPLS based network adding short-term capacity over Internet using IPsec, there are secure connection between edge nodes and RR, via private path, TLS, DTLS, etc. The authentication of peer nodes is managed by the RR. More importantly, when an edge node needs to establish multiple IPsec tunnels to many different edge nodes, all the management information can be multiplexed into the secure management tunnel between RR and the edge node. Therefore, there is reduced amount of authentication in a BGP Controlled SDWAN network.

Client VPNs are configured via VRFs, just like the configuration of the existing MPLS VPN. The IPsec equivalent traffic selectors for

local and remote routes is achieved by importing/exporting VPN Route Targets. The binding of client routes to IPsec SA is dictated by policies. As the result, the IPsec configuration for a BGP controlled SDWAN (with mixed MPLS VPN) can be simplified as the following:

- SDWAN controller has authority to authenticate edges and peers. Remote Peer association is controlled by the SDWAN Controller (RR)
- The IKEv2 proposals including the IPsec Transform set can be sent directly to Peer or incorporated with BGP UPDATE.
- BGP UPDATE: Announce the client route reachability for all permitted parallel tunnels/paths.
  - o No need to run multiple routing protocols in each IPsec tunnel.
- using importing/exporting Route Targets under each client VPN (VRF) to achieve the traffic selection (or permission) among clients' routes attached to multiple edge nodes.

### **3.3. Client UPDATE and Hybrid Underlay Tunnel UPDATE**

As described in [[SDWAN-BGP-USAGE](#)], two separate BGP UPDATE messages are used for SDWAN Edge Discovery:

- UPDATE U1 for the attached client routes,  
This UPDATE is for a SDWAN node to advertise the attached client routes to remote peers. This UPDATE will continue using the existing BGP AFI/SAFI for IP or VPN. Detailed underlay tunnel specification can be recursively resolved by using the Recursive Next Hop Resolution as specified by the section 8 of [[Tunnel-Encap](#)].

A new Tunnel Type (SDWAN-Hybrid) needs to be added, to be used by Encapsulation Extended Community or the Tunnel-Encap Path Attribute [Tunnel-encap] to indicate mixed underlay networks.

- UPDATE U2, advertised by the Next hop address of the UPDATE U1 to propagate the properties of the tunnels terminated at the edge node.  
This UPDATE is for an edge node to advertise the properties of directly attached underlay networks, including the NAT information, pre-configured IPsec SA identifiers, and/or the underlay network ISP information. This UPDATE can also include

the detailed IPsec SA attributes, such as keys, nonce, encryption algorithms, etc.

The UPDATE U2 is for peers to discover remote node's underlay network properties.

In the following figure: there are four types underlay paths between C-PE1 and C-PE2:

- a) MPLS-in-GRE path.
- b) node-based IPsec tunnel [2.2.2.2<->1.1.1.1].
- c) port-based IPsec tunnel [192.0.0.1 <-> 192.10.0.10]; and
- d) port-based IPsec tunnel [172.0.0.1 <-> 160.0.0.1].

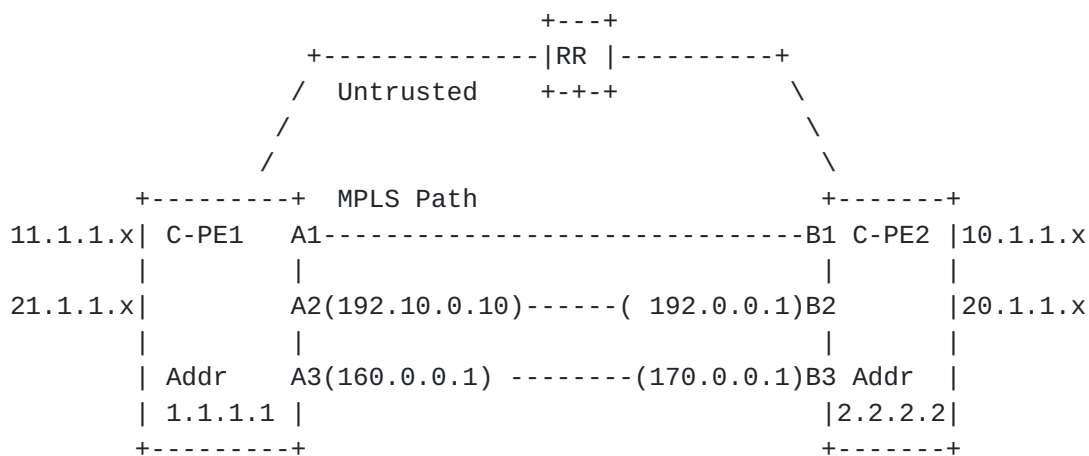


Figure 1: Hybrid SDWAN

C-PE2 uses UPDATE U1 to advertise the attached client routes:

UPDATE U1:

Extended community: RT for SDWAN VPN 1  
 NLRI: AFI=? & SAFI = 1/1  
 Prefix: 10.1.1.x; 20.1.1.x

NextHop: 2.2.2.2 (C-PE2)  
Encapsulation Extended Community: tunnel-type=SDWAN-hybrid  
Color Extended Community: RED

The UPDATE U1 is recursively resolved to the UPDATE U2 which specifies the detailed hybrid WAN underlay Tunnels terminated at the C-PE2:

UPDATE U2:

NLRI: SAFI = SDWAN  
(With Color RED encoded in the NLRI Site-Property field)  
Prefix: 2.2.2.2  
Tunnel encapsulation Path Attribute [type=SDWAN-Hybrid]  
  IPSec SA for 192.0.0.1  
  Tunnel-End-Point Sub-TLV [[Section 3.1](#) of Tunnel-encap]  
  IPsec SA sub-TLV [See the [Section 5](#)]  
Tunnel encapsulation Path Attribute [type=SDWAN-Hybrid]  
  IPSec SA for 170.0.0.1  
  Tunnel-End-Point Sub-TLV /\*the address\*/  
  IPsec SA sub-TLV  
  
Tunnel Encap Attr MPLS-in-GRE [type=SDWAN-Hybrid]  
  Sub-TLV for MPLS-in-GRE [[Section 3.2.6](#) of Tunnel-encap]

Note: [[Tunnel-Encap](#)] [Section 11](#) specifies that each Tunnel Encap Attribute can only have one Tunnel-End-Point sub-TLV. Therefore, two separate Tunnel Encap Attributes are needed to indicate that the client routes can be carried by either one.

### **3.4. Edge Node Discovery**

The basic scheme of SDWAN Edge node discovery using BGP consists of:

- Secure connection to a SDWAN controller (i.e. RR in this context):  
For a SDWAN edge with both MPLS and IPsec path, the edge node should already have secure connection to its controller, i.e. RR in this context. For a remote SDWAN edge that is only accessible via Internet, the SDWAN edge node, upon power up, establishes a secure tunnel (such as TLS, SSL) with the SDWAN central controller whose address is preconfigured on the edge node. The central controller will inform the edge node of its



local RR. The edge node will establish a transport layer secure session with the RR (such as TLS, SSL).

- The Edge node will advertise its own properties to its designated RR via the secure connection.
- The RR propagates the received information to the authorized peers.
- The authorized peers can establish the secure data channels (IPsec) and exchange more information among each other.

For a SDWAN deployment with multiple RRs, it is assumed that there are secure connections among those RRs. How secure connections being established among those RRs is the out of the scope of the current draft. The existing BGP UPDATE propagation mechanisms control the edge properties propagation among the RRs.

For some special environment where the communication to RR are highly secured, [SDN-IPsec] IKE-less can be deployed to simplify IPsec SA establishment among edge nodes.

#### **4. BGP UPDATE to Support SDWAN Segmentation**

##### **4.1. SDWAN Segmentation, SDWAN Virtual Topology and Client VPN**

In SDWAN deployment, "SDWAN Segmentation" is a frequently used term, referring to partitioning a network to multiple sub-networks, just like what MPLS VPN does. "SDWAN Segmentation" is achieved by creating SDWAN virtual topologies and SDWAN VPNs. A SDWAN virtual topology consists of a set of edge nodes and the tunnels, including both IPsec tunnels and/or MPLS VPN tunnels), interconnecting those edge nodes.

A SDWAN VPN is same as a client VPN, which is configured in the same way as the VRFs on PEs of a MPLS VPN. One SDWAN client VPN can be mapped to one or multiple SD-WAN virtual topologies. How a Client VPN is mapped to a SDWAN virtual topology is governed by policies from the SDWAN controller.

Each SDWAN edge node may need to support multiple VPNs. Just like Route Target is used to distinguish different MPLS VPNs, SDWAN VPN ID is used to differentiate the SDWAN VPNs. For example, in the picture below, the "Payment-Flow" on C-PE2 is only mapped to the

virtual topology of C-PEs to/from Payment Gateway, whereas other flows can be mapped to a multipoint to multipoint virtual topology.

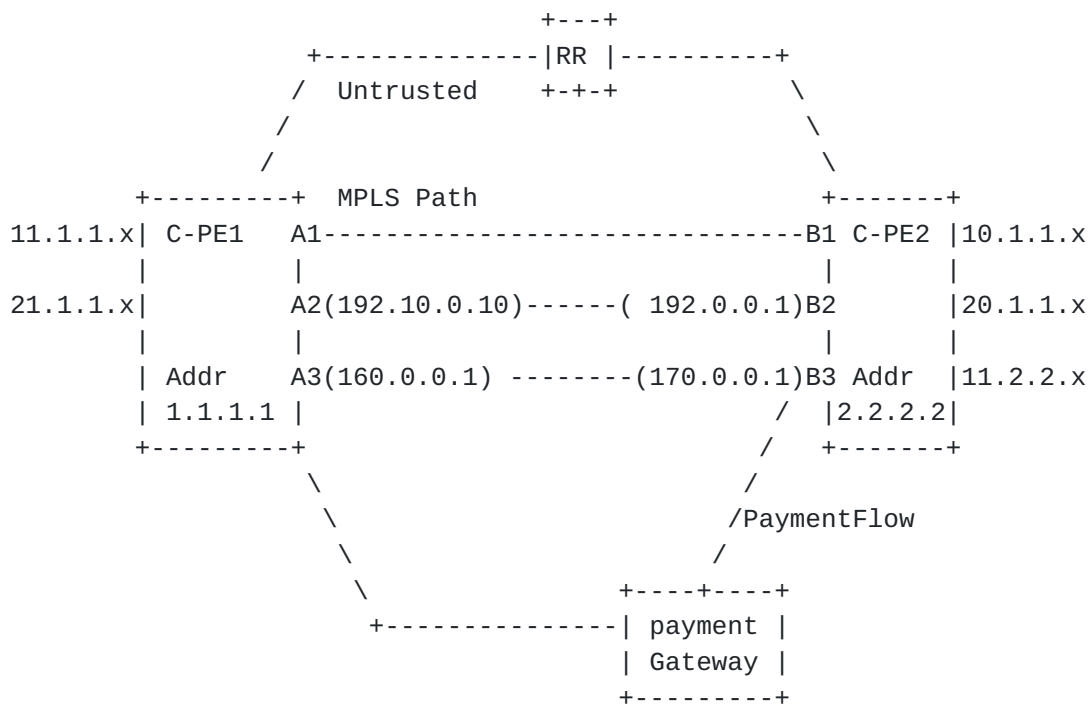


Figure 2: SDWAN Virtual Topology & VPN

#### 4.2. Constrained Propagation of Edge Capability

BGP has built-in mechanism to dynamically achieve the constrained distribution of edge information. [RFC4684](#) describes the BGP RT constrained distribution. In a nutshell, a SDWAN edge sends RT Constraint (RTC) NLRI to the RR for the RR to install an outbound route filter, as shown in the figure below:

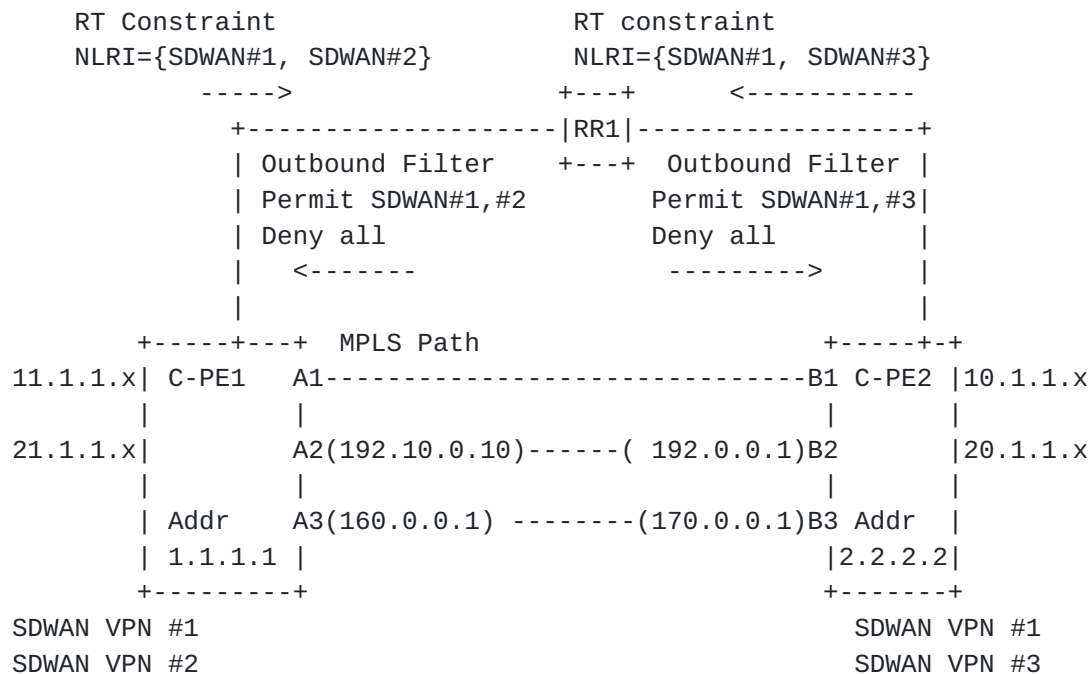


Figure 3: Constraint propagation of Edge Property

However, a SDWAN overlay network can span across untrusted networks, RR can't trust the RT Constraint (RTC) NLRI BGP UPDATE from any nodes. RR can only process the RTC NLRI from authorized peers for a SDWAN VPN.

It is out of the scope of this document on how RR is configured with the policies to filter out unauthorized nodes for specific SDWAN VPNs.

When the RR receives BGP UPDATE from an edge node, it propagates the received UPDATE message to the nodes that are in the Outbound Route filter for the specific SDWAN VPN.

#### 4.3. SDWAN VPN ID in BGP Update

SDWAN VPN is represented by the SDWAN VPN ID in the BGP Extended Community.

A different Extended Community Type is used to represent SDWAN VPN ID, so that routes that can only be carried by MPLS path can be differentiated from routes that can be carried by hybrid paths.

Encoding:

[RFC4360](#): Extended Community for SDWAN VPN ID:

```

  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 high      | Type low(*)   |                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Value                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

  0 1 2 3 4 5 6 7
+---+---+---+---+---+
|I|T| 6-bit val |
+---+---+---+---+---+

```

The high-order octet of the Type Field

T bit =0 (transitive) when SDWAN edge sends to its RR which then propagates to remote peers based on outbound filters.

[RFC4760](#) states that Route Target community is transitive

For SDWAN, an edge receiving the SDWAN Update shouldn't forward it to other nodes.

T bit =1 (non-transitive) when RR propagates the UPDATE to SDWAN EDGE

[IANA Consideration:

Following the encoding scheme specified by [RFC7153](#), need IANA to assign the following values for the "Type High" Octet:

- Transitive (when edge announce the advertisement to its RR): 0x0A, which is the number after 0x08 for Flow Spec Redirect.
- Non Transitive (when RR send to remote edges): 0x4A

Request a new value of the low-order octet of the Type field for this community (different from the VPN Route Target 0x02).

]

#### 4.4. SDWAN VPN ID in Data Plane

From data plane perspective, packets from different SDWAN VPNs need to have their corresponding SDWAN VPN identifier encoded in the header.

For a SDWAN edge node which can be reached by both MPLS and IPsec paths, the client packets reached by MPLS network will be encoded with the MPLS Labels based on the scheme specified by [RFC8277](#).

For GRE Encapsulation within IPsec tunnel, the GRE key field can be used to carry the SDWAN VPN ID. For NVO (VxLAN, GENEVE, etc.) encapsulation within the IPsec tunnel, Virtual Network Identifier (VNI) field is used to carry the SDWAN VPN ID.

### 5. Hybrid Underlay Tunnel UPDATE

The hybrid underlay tunnel UPDATE is to advertise the detailed properties of hybrid types of tunnels terminated at a SDWAN edge node.

A client route UPDATE is recursively tied to an underlay tunnel UPDATE by the Color Extended Community included in client route UPDATE.

#### 5.1. NLRI for Hybrid Underlay Tunnel Update

A new NLRI is introduced within the MP\_REACH\_NLRI Path Attribute of [RFC4760](#), for advertising the detailed properties of hybrid types of tunnels terminated at the edge node, with SAFI=SDWAN (code = 74):

```
+-----+
|  NLRI Length   | 1 octet
+-----+
|  Site-Type     | 1 Octet
+-----+
|  Port-Local-ID | 4 octets
+-----+
|  SDWAN-Color   | 4 octets
+-----+
|  SDWAN-Node-ID | 4 or 16 octets
+-----+
```

where:

- NLRI Length: 1 octet of length expressed in bits as defined in [\[RFC4760\]](#).
- Site Type: 1 octet value. The SDWAN Site Type defines the different types of Site IDs to be used in the deployment. The draft defines the following types:
  - Site-Type = 1: For simple deployment, such as all edge nodes under one SDWAN management system, a simple identifier is enough for the SDWAN management to map the site to its precise geolocation.
  - Site-Type = 2: to indicate that the value in the site-ID is locally significant, therefore, need a Geo-Loc Sub-TLV to fully describe the accurate location of the node. This is for large SDWAN heterogeneous deployment where Site IDs has to be described by proper Geo-location of the Edge Nodes [\[LISP-GEOLoc\]](#).
- Port local ID: SDWAN edge node Port identifier, which can be locally significant. The detailed properties about the network connected to the port are further encoded in the Tunnel Path Attribute. If the SDWAN NLRI applies to multiple ports, this field is NULL.
- SDWAN-Color: is used to correlate with the Color-Extended-community included in the client routes UPDATE.
- SDWAN Edge Node ID: a routable address (IPv4 or IPv6) within the WAN to reach this node or port.

[Editor's note on using SDWAN SAFI for the underlay network property advertisement:

SDWAN SAFI [\[IANA assigned =74\]](#) is used instead of IP SAFI in the MP-NLRI [\[RFC4760\]](#) Path Attribute to advertise the underlay network properties to emphasize that the address in the NLRI is NOT client addresses.

If the same IP SAFI used, receiver needs to add extra logic to differentiate regular BGP MP-NLRI client routes advertisement from the SDWAN underlay network properties



IPsec-SA-ID Sub-TLV: the length of the sub-TLV is 4 octets. The following is the structure of the Value field of the IPsec-SA-ID sub-TLV.

```

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
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                               IPsec SA Identifier                               |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

The IPsec SA identifier (4 Octet) is for cross reference the IPsec SA attributes being pre-configured or advertised by another UPDATE [[Section 7](#)].

If the client traffic needs to be encapsulated in a specific type within the IPsec ESP Tunnel, such as GRE or VxLAN, etc., the corresponding Tunnel-Encap Sub-TLV needs to be appended right after the IPsec-SA-ID Sub-TLV.

IPsec-SA-Group Sub-TLV is for the scenario that multiple IPsec SAs have the same inner encapsulation. Multiple IPsec SA IDs are included in the IPsec-ID-Group Sub-TLV. If different inner encapsulation is desired within IPsec tunnels, multiple IPsec-SA-Group Sub-TLVs can be included within one Tunnel Encap Path Attribute. The following is the structure of the Value field of the IPsec-SA-Group sub-TLV.

```

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
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|      InnerEncapType      |      reserved      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|      IPsec SA Identifier #1      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
~      IPsec SA Identifier #n      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
IPsec-SA-Group Sub-TLV

```

InnerEncapType (2 octet) indicates the encapsulation type for the payload within the IPsec ESP Tunnel. The Inner Encap Type value will take the value specified by the IANA Consideration Section (12.5) of [[Tunnel-Encap](#)]:

- types 8 (VXLAN), 9 (NVGRE), 11 (MPLS-in-GRE), and 12 (VXLAN-GPE) in the "BGP Tunnel Encapsulation Tunnel Types" registry.



- types 1 (L2TPv3), 2 (GRE), and 7 (IP in IP) in the "BGP Tunnel Encapsulation Tunnel Types" registry.

For each of the Tunnel Types specified, the detailed encapsulation value field as specified by Section 3.2 of [[Tunnel-Encap](#)] is appended right after the IPsec Sub-TLV.

The Tunnel Ending Point Sub-TLV specified by the Section 3.1 of [[Tunnel-Encap](#)] has to be attached to identify the IPsec Tunnel terminating address.

There can be multiple IPsec tunnels terminating at one WAN port or at one node, e.g. one tunnel for going to destination "A", another one for going to destination "B". Use SDWAN for retail industry as an example, it is necessary for all shops at any location to only exchange Payment System traffic with the Payment Gateway, while other traffic can be exchanged with any nodes.

Therefore, there could be multiple IPsec Sub-TLVs bound with one Tunnel Ending Point Sub-TLV.

#### **5.3.1. Encoding example #1 of using IPsec-SA-ID Sub-TLV**

This section provides an encoding example for the following scenario:

- there are three IPsec SAs terminated at the same WAN Port address (or the same node address)
- Two of the IPsec SAs use GRE (value =2) as Inner Encapsulation within the IPsec Tunnel
- One of the IPsec SA uses VxLAN (value = 8) as the Inner Encapsulation within its IPsec Tunnel.

Here is the encoding for the scenario:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Tunnel-Type =SDWAN-Hybrid          |          Length =          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          Tunnel-end-Point Sub-TLV          |          |
~                                          ~

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| subTLV-Type                        = IPsec-SA-group      |          Length
=          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| IPsec-SA-group:InnerEncapType=2|          Reserved          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          IPsec SA Identifier = 1          |          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          IPsec SA Identifier = 2          |          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          GRE-KEY (4 Octets)          |          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|subTLV-Type = IPsec-SA-ID          |          Length= 4          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          IPsec SA Identifier = 1          |          |
+-----+-----+-----+-----+-----+-----+-----+-----+
~          VxLAN Sub-TLV          ~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

The Length of the Tunnel-Type = SDDWAN-Hybrid is the sum of the following:

- Tunnel-end-point sub-TLV total length
- the IPsec-SA-Group Sub-TLV length + 4 (the two octets for InnerEncapType + the two octets for the Length field)
- GRE-Key Length (4)
- The IPsec-SA-ID Sub-TLV length: 4
- The VxLAN sub-TLV total length

### 5.3.2. Encoding Example #2 of using IPsec-SA-ID

If IPsec SAs are terminating at different addresses, then multiple Tunnel Encap Attributes have to be included.

The encoding example for the Figure 1:

- there is one IPsec SA terminating at the WAN Port address 192.0.0.1; and another IPsec SA terminating at WAN Port 170.0.0.1;

- Both IPsec SAs use GRE (value =2) as Inner Encapsulation within the IPsec Tunnel

```

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
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Tunnel-Type =SDWAN-Hybrid      |      Length =      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|      Tunnel-end-Point Sub-TLV      |      |
~      for 192.0.0.1      ~
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
~      IPsec-SA-ID sub-TLV #1      ~
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
~      GRE Sub-TLV      ~
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Tunnel-Type =SDWAN-Hybrid      |      Length =      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|      Tunnel-end-Point Sub-TLV      |      |
~      for 170.0.0.1      ~
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
~      IPsec-SA-ID sub-TLV #2      ~
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
~      GRE sub-TLV      ~
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

#### 5.4. Extended Port Sub-TLV

When a SDWAN edge node is connected to an underlay network via a port behind NAT devices, traditional IPsec uses IKE for NAT negotiation. The location of a NAT device can be such that:

- Only the initiator is behind a NAT device. Multiple initiators can be behind separate NAT devices. Initiators can also connect to the responder through multiple NAT devices.
- Only the responder is behind a NAT device.
- Both the initiator and the responder are behind a NAT device.

The initiator's address and/or responder's address can be dynamically assigned by an ISP or when their connection crosses a

dynamic NAT device that allocates addresses from a dynamic address pool.

Because one SDWAN edge can connect to multiple peers via one underlay network, the pair-wise NAT exchange as IPsec's IKE is not efficient. In BGP Controlled SDWAN, NAT information of a WAN port is advertised to its RR in the BGP UPDATE message. It is encoded as an Extended sub-TLV that describes the NAT property if the port is behind a NAT device.

A SDWAN edge node can inquire STUN (Session Traversal of UDP Through Network Address Translation [RFC 3489](https://tools.ietf.org/html/rfc3489)) Server to get the NAT property, the public IP address and the Public Port number to pass to peers.

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|Port Ext Type | EncapExt subTLV Length          |I|O|R|R|R|R|R|
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| NAT Type     | Encap-Type |Trans networkID|      RD ID      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|
|          Local IP Address
|          32-bits for IPv4, 128-bits for Ipv6
|          ~~~~~~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          Local Port
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          Public IP
|          32-bits for IPv4, 128-bits for Ipv6
|          ~~~~~~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          Public Port
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          ISP-Sub-TLV
~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Where:

- o Port Ext Type: indicate it is the Port Ext SubTLV.
- o PortExt subTLV Length: the length of the subTLV.

- o Flags:

- I bit (CPE port address or Inner address scheme)
  - If set to 0, indicate the inner (private) address is IPv4.
  - If set to 1, it indicates the inner address is IPv6.
- O bit (Outer address scheme):
  - If set to 0, indicate the public (outer) address is IPv4.
  - If set to 1, it indicates the public (outer) address is IPv6.
- R bits: reserved for future use. Must be set to 0 now.

- o NAT Type.without NAT; 1:1 static NAT; Full Cone; Restricted Cone; Port Restricted Cone; Symmetric; or Unknown (i.e. no response from the STUN server).
- o Encap Type.the supported encapsulation types for the port facing public network, such as IPsec+GRE, IPsec+VxLAN, IPsec without GRE, GRE (when packets don't need encryption)
- o Transport Network ID.Central Controller assign a global unique ID to each transport network.
- o RD ID.Routing Domain ID.need to be global unique.
- o Local IP.The local (or private) IP address of the port.
- o Local Port.used by Remote SDWAN edge node for establishing IPsec to this specific port.
- o Public IP.The IP address after the NAT. If NAT is not used, this field is set to NULL.
- o Public Port.The Port after the NAT. If NAT is not used, this field is set to NULL.

#### 5.5. ISP of the Underlay network Sub-TLV

The purpose of the Underlay network Sub-TLV is to carry the ISP WAN port properties with SDWAN SAFI NLRI. It would be treated as optional Sub-TLV. The BGP originator decides whether to include this Sub-TLV along with the SDWAN NLRI. If this Sub-TLV is present, it would be processed by the BGP receiver and to determine what local policies to apply for the remote end point of the Underlay tunnel.

The format of this Sub-TLV is as follows:

```

  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   |      Flag      |      Reserved   |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|Connection Type|   Port Type   |      Port Speed      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

Where:

Type: To be assigned by IANA

Length: 6 bytes.

Flag: a 1 octet value.

Reserved: 1 octet of reserved bits. It SHOULD be set to zero on transmission and MUST be ignored on receipt.

Connection Type: There are two different types of WAN Connectivity. They are listed below as:

```

Wired - 1
WIFI - 2
LTE - 3
5G - 4

```

Port Type: There are different types of ports. They are listed Below as:

```

Ethernet - 1
Fiber Cable - 2
Coax Cable - 3
Cellular - 4

```

Port Speed: The port seed is defined as 2 octet value. The values are defined as Gigabit speed.

## 6. IPsec SA Property Sub-TLVs

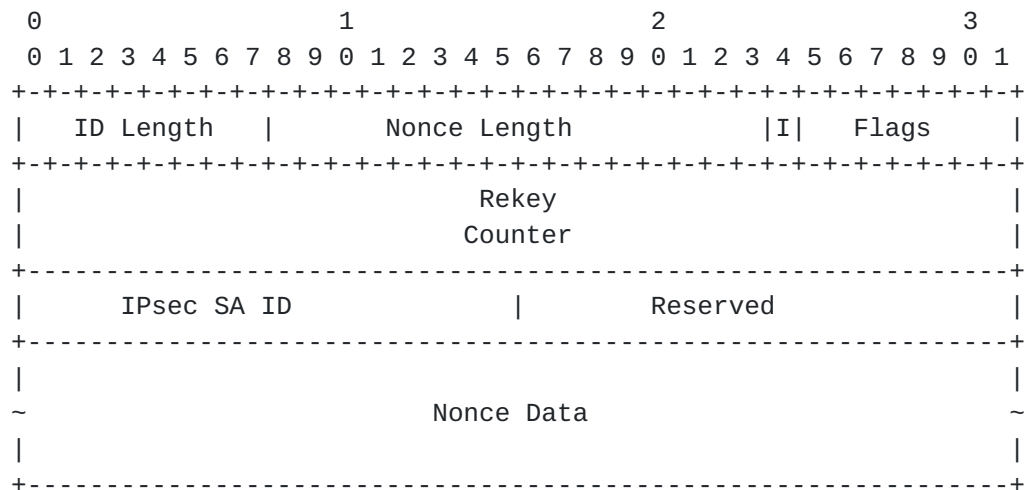
This section describes the detailed IPsec SA properties sub-TLVs.

### 6.1. IPsec SA Nonce Sub-TLV

The Nonce Sub-TLV is based on the Base DIM sub-TLV as described the Section 6.1 of [\[SECURE-EVPN\]](#). IPsec SA ID is added to the sub-TLV, which is to be referenced by the client route NLRI Tunnel Encap Path Attribute for the IPsec SA. The following fields are removed because:

- the Originator ID is carried by the NLRI,
- the Tenant ID is represented by the SDWAN VPN ID Extended Community, and
- the Subnet ID are carried by the BGP route UPDATE.

The format of this Sub-TLV is as follows:

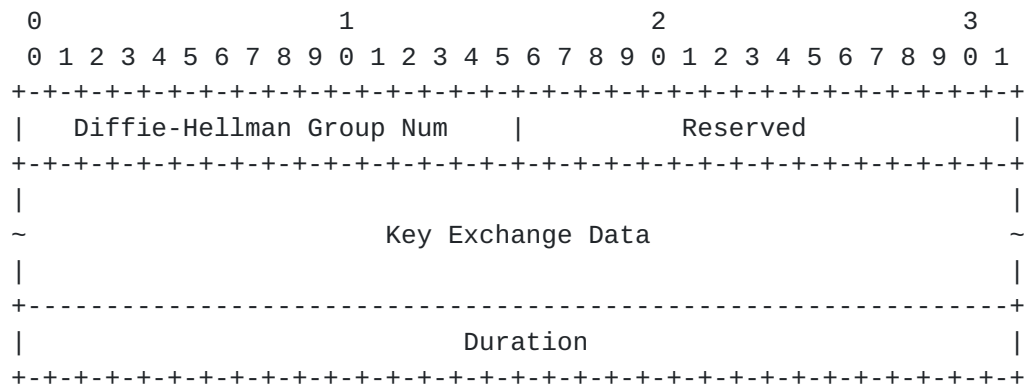


IPsec SA ID - The 2 bytes IPsec SA ID could 0 or non-zero values. It is cross referenced by client route's IPsec Tunnel Encap IPsec-SA-ID or IPsec-SA-Group Sub-TLV in [Section 5](#). When there are multiple IPsec SAs terminated at one address, such as WAN port address or the node address, they are differentiated by the different IPsec SA IDs.

### 6.2. IPsec Public Key Sub-TLV

The IPsec Public Key Sub-TLV is derived from the Key Exchange Sub-TLV described in [\[SECURE-EVPN\]](#) with an addition of Duration field to define the IPsec SA life span. The edge nodes would pick the shortest duration value between the SDWAN SAFI pairs.

The format of this Sub-TLV is as follows:



### 6.3. IPsec SA Proposal Sub-TLV

The IPsec SA Proposal Sub-TLV is to indicate the number of Transform Sub-TLVs. This Sub-TLV aligns with the sub-TLV structure from [\[SECURE-VPN\]](#)

The Transform Sub-sub-TLV will follow the [section 3.3.2 of RFC7296](#).

### 6.4. Simplified IPsec Security Association sub-TLV

For a simple SDWAN network with edge nodes supporting only a few pre-defined encryption algorithms, a simple IPsec sub-TLV can be used to encode the pre-defined algorithms, as below:



```

 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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|IPsec-simType |IPsecSA Length          | Flag          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Transform    | Mode                  | AH algorithms |ESP algorithms |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|              ReKey Counter (SPI)              |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| key1 length  |          Public Key                  ~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| key2 length  |          Nonce                      ~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|              Duration                  |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Where:

- o IPsec-SimType: The type value has to be between 128~255 because IPsec-SA subTLV needs 2 bytes for length to carry the needed information.
- o IPsec-SA subTLV Length (2 Byte): 25 (or more)
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o Transform (1 Byte): the value can be AH, ESP, or AH+ESP.
- o IPsec Mode (1 byte): the value can be Tunnel Mode or Transport mode
- o AH algorithms (1 byte): AH authentication algorithms supported, which can be md5 | sha1 | sha2-256 | sha2-384 | sha2-512 | sm3. Each SDWAN edge node can have multiple authentication algorithms; send to its peers to negotiate the strongest one.
- o ESP (1 byte): ESP authentication algorithms supported, which can be md5 | sha1 | sha2-256 | sha2-384 | sha2-512 | sm3. Each SDWAN edge node can have multiple authentication algorithms; send to its peers to negotiate the strongest one. Default algorithm is AES-256.
  - o When node supports multiple authentication algorithms, the initial UPDATE needs to include the "Transform Sub-TLV" described by [\[SECURE-EVPN\]](#) to describe all of the algorithms supported by the node.

- o Rekey Counter (Security Parameter Index): 4 bytes
- o Public Key: IPsec public key
- o Nonce: IPsec Nonce
- o Duration: SA life span.

### 6.5. IPsec SA Encoding Examples

For the Figure 1 in [Section 3](#), C-PE2 needs to advertise its IPsec SA associated attributes, such as the public keys, the nonce, the supported encryption algorithms for the IPsec tunnels terminated at 192.0.0.1, 170.1.1.1 and 2.2.2.2 respectively.

Using the IPsec Tunnel [ISP4: 160.0.0.1 <-> ISP2:170.0.0.1] as an example: C-PE1 needs to advertise the following attributes for establishing the IPsec SA:

SDWAN Node ID

SDWAN Color

Tunnel Encap Attr (Type=SDWAN-Hybrid)

Extended Port Sub-TLV for information about the Port  
(including ISP Sub-TLV for information about the ISP2)

IPsec SA Nonce Sub-TLV,

IPsec SA Public Key Sub-TLV,

IPsec SA Sub-TLV for the supported transforms

{Transforms Sub-TLV - Trans 2,  
Transforms Sub-TLV - Trans 3}

C-PE2 needs to advertise the following attributes for establishing IPsec SA:

SDWAN Node ID

SDWAN Color

Tunnel Encap Attr (Type=SDWAN-Hybrid)

Extended Port Sub-TLV (including ISP Sub-TLV for information  
about the ISP2)

IPsec SA Nonce Sub-TLV,

IPsec SA Public Key Sub-TLV,

IPsec SA Sub-TLV for the supported transforms

{Transforms Sub-TLV - Trans 2,  
Transforms Sub-TLV - Trans 4}

As both end points support Transform #2, the Transform #2 will be used for the IPsec Tunnel [ISP4: 160.0.0.1 <-> ISP2:170.0.0.1].

## 7. Error & Mismatch Handling

Each C-PE device advertises SDWAN SAFI Underlay NLRI to the other C-PE devices via BGP Route Reflector to establish pairwise SAs between itself and every other remote C-PEs. During the SAFI NLRI advertisement, the BGP originator would include either simple IPsec Security Association properties defined in IPsec SA Sub-TLV based on IPsec-SA-Type = 1 or full-set of IPsec Sub-TLVs including Nonce, Public Key, Proposal and number of Transform Sub-TLVs based on IPsec-SA-Type = 2.

The C-PE devices would compare the IPsec SA attributes between the local and remote WAN ports. If there is a match on the SA Attributes between the two ports, the IPsec Tunnel would be established.

The C-PE devices would not try to negotiate the base IPsec-SA parameters between the local and the remote ports in the case of simple IPsec SA exchange or the Transform sets between local and remote ports if there is a mismatch on the Transform sets in the case of full-set of IPsec SA Sub-TLVs.

As an example, using the Figure 1 in [Section 3](#), to establish IPsec Tunnel between C-PE1 and C-PE2 WAN Ports A2 and B2 [A2: 192.10.0.10 <-> B2:192.0.0.1]:

C-PE1 needs to advertise the following attributes for establishing the IPsec SA:

- NH: 192.10.0.10

- SDWAN Node ID

- SDWAN-Site-ID

- Tunnel Encap Attr (Type=SDWAN)

  - ISP Sub-TLV for information about the ISP3

  - IPsec SA Nonce Sub-TLV,

  - IPsec SA Public Key Sub-TLV,

  - Proposal Sub-TLV with Num Transforms = 1

    - {Transforms Sub-TLV - Trans 1}

C-PE2 needs to advertise the following attributes for establishing IPsec SA:

- NH: 192.0.0.1
- SDWAN Node ID
- SDWAN-Site-ID
- Tunnel Encap Attr (Type=SDWAN)
  - ISP Sub-TLV for information about the ISP1
  - IPsec SA Nonce Sub-TLV,
  - IPsec SA Public Key Sub-TLV,
  - Proposal Sub-TLV with Num Transforms = 1
  - {Transforms Sub-TLV - Trans 2}

As there is no matching transform between the WAN ports A2 and B2 in C-PE1 and C-PE2 respectively, there will be no IPsec Tunnel be established.

## **8. Manageability Considerations**

TBD - this needs to be filled out before publishing

## **9. Security Considerations**

The document describes the encoding for SDWAN edge nodes to advertise its properties to their peers to its RR, which propagates to the intended peers via untrusted networks.

The secure propagation is achieved by secure channels, such as TLS, SSL, or IPsec, between the SDWAN edge nodes and the local controller RR.

[More details need to be filled in here]

## **10. IANA Considerations**

This document requires the following IANA actions.

- o Hybrid (SDWAN) Overlay SAFI = 74 assigned by IANA
- o SDWAN VPN ID Extended Community type
- o IPsec-SA-ID Sub-TLV Type

- o IPsec-SA-Group Sub-TLV Type

## **11. References**

### 11.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

### 11.2. Informative References

- [RFC8192] S. Hares, et al, "Interface to Network Security Functions (I2NSF) Problem Statement and Use Cases", July 2017
- [RFC5521] P. Mohapatra, E. Rosen, "The BGP Encapsulation Subsequent Address Family Identifier (SAFI) and the BGP Tunnel Encapsulation Attribute", April 2009.
- [CONTROLLER-IKE] D. Carrel, et al, "IPsec Key Exchange using a Controller", [draft-carrel-ipsecme-controller-ike-01](#), work-in-progress.
- [LISP-GEOLoc] D. Farinacci, "LISP Geo-Coordinate Use-Case", [draft-farinacci-lisp-geo-09](#), April 2020.
- [SDN-IPSEC] R. Lopez, G. Millan, "SDN-based IPsec Flow Protection", [draft-ietf-i2nsf-sdn-ipsec-flow-protection-07](#), Aug 2019.
- [SECURE-EVPN] A. Sajassi, et al, "Secure EVPN", [draft-sajassi-bess-secure-evpn-02](#), July 2019.
- [[Tunnel-Encap](#)] E. Rosen, et al, "The BGP Tunnel Encapsulation Attribute", [draft-ietf-idr-tunnel-encaps-09](#), Feb 2018.
- [VPN-over-Internet] E. Rosen, "Provide Secure Layer L3VPNs over Public Infrastructure", [draft-rosen-bess-secure-l3vpn-00](#), work-in-progress, July 2018

[DMVPN] Dynamic Multi-point VPN:

<https://www.cisco.com/c/en/us/products/security/dynamic-multipoint-vpn-dmvpn/index.html>

[DSVPN] Dynamic Smart VPN:

<http://forum.huawei.com/enterprise/en/thread-390771-1-1.html>

[ITU-T-X1036] ITU-T Recommendation X.1036, "Framework for creation, storage, distribution and enforcement of policies for network security", Nov 2007.

[Net2Cloud-Problem] L. Dunbar and A. Malis, "Seamless Interconnect Underlay to Cloud Overlay Problem Statement", [draft-dm-net2cloud-problem-statement-02](#), June 2018

[Net2Cloud-gap] L. Dunbar, A. Malis, and C. Jacquenet, "Gap Analysis of Interconnecting Underlay with Cloud Overlay", [draft-dm-net2cloud-gap-analysis-02](#), work-in-progress, Aug 2018.

[Tunnel-Encap] E. Rosen, et al "The BGP Tunnel Encapsulation Attribute", [draft-ietf-idr-tunnel-encaps-10](#), Aug 2018.

## **12. Acknowledgments**

Acknowledgements to Wang Haibo, Hao Weiguo, and ShengCheng for implementation contribution; Many thanks to Jim Guichard, John Scudder, and Donald Eastlake for their review and contributions.

This document was prepared using 2-Word-v2.0.template.dot.



Authors' Addresses

Linda Dunbar  
Futurewei  
Email: ldunbar@futurewei.com

Sue Hares  
Hickory Hill Consulting  
Email: shares@ndzh.com

Robert Raszuk  
Email: robert@raszuk.net

Kausik Majumdar  
CommScope  
Email: Kausik.Majumdar@commscope.com