

July 2, 2018

Gap Analysis of Interconnecting Underlay with Cloud Overlay
[draft-dm-net2cloud-gap-analysis-00](#)

Abstract

This document analyzes the technological gaps when using SD-WAN to interconnect workloads & apps hosted in various locations, especially cloud data centers when the network service providers do not have or have limited physical infrastructure to reach the locations [Net2Cloud-problem].

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#). This document may not be modified, and derivative works of it may not be created, except to publish it as an RFC and to translate it into languages other than English.

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/lid-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 December 2, 2018.

Copyright Notice

Copyright (c) 2018 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1.	Introduction.....	3
2.	Conventions used in this document.....	3
3.	Gap Analysis of CPEs Registration Protocol.....	4
4.	Gap Analysis in aggregating VPN paths and Internet paths.....	4
4.1.	Gap analysis of Using BGP to cover SD-WAN paths.....	6
4.2.	Gaps in preventing attacks to CPEs from their Internet ports	7
5.	Gap analysis of CPEs not directly connected to VPN PEs.....	8
5.1.	Gap Analysis of Floating PEs to connect to Remote CPEs...	10
5.2.	NAT Traversing.....	10
5.3.	Complication of use BGP between PE and remote CPEs via Internet.....	10
5.4.	Designated Forwarder to the remote edges.....	11
5.5.	Traffic Path Management.....	12
6.	Manageability Considerations.....	12
7.	Security Considerations.....	12
8.	IANA Considerations.....	12
9.	References.....	12
9.1.	Normative References.....	13
9.2.	Informative References.....	13

10.	Acknowledgments.....	14
---------------------	----------------------	--------------------

[1.](#) Introduction

[Net2Cloud-Problem] describes the problems of enterprises face today in transitioning their IT infrastructure to support digital economy, such as connecting enterprises' branch offices to dynamic workloads in Cloud DCs.

This document analyzes the technological gaps to interconnect dynamic workloads & apps hosted in various locations, especially in cloud data centers to which the network service providers may not have or have limited physical infrastructure to reach.

[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 SD-WAN controller to manage SD-WAN overlay path creation/deletion and monitoring the path conditions between two sites.

CPE-Based VPN: Virtual Private Secure network formed among CPEs. This is to differentiate from most commonly used PE based VPNs

OnPrem: On Premises data centers and branch offices

SD-WAN: Software Defined Wide Area Network, which can mean many different things. In this document, "SD-WAN" refers to the solutions specified by ONUG (Open Network User Group), which build point-to-point IPsec overlay paths between two end-points (or branch offices) that need to intercommunicate.

3. Gap Analysis of CPEs Registration Protocol

SD-WAN, conceived in ONUG (Open Network User Group) a few years ago as way to aggregate multiple connections between any two points, has emerged as an on-demand technology to securely interconnect the OnPrem branches with the workloads instantiated in Cloud DCs that do not have MPLS VPN PE co-located or have very limited bandwidths.

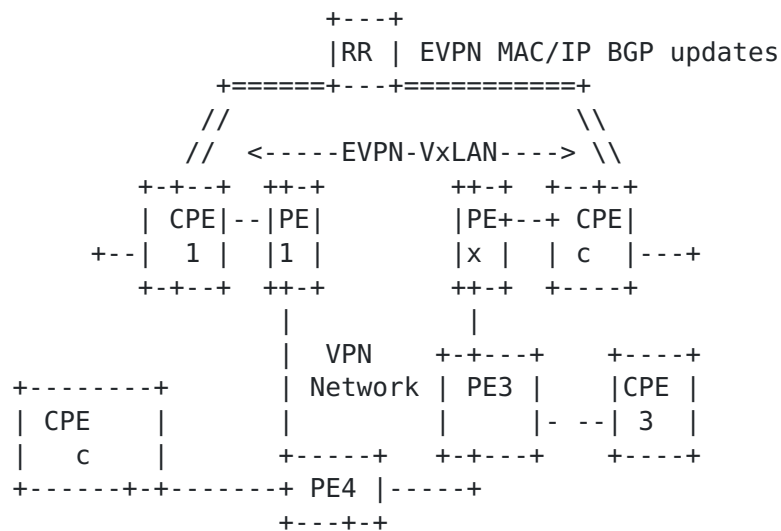
Some SD-WAN networks use the NHRP protocol [[RFC2332](#)] to register SD-WAN endpoints with a "Controller" (or NHRP server), which then has the ability to map a private VPN address to a public IP address of the destination node. DSVPN [[DSVPN](#)] or DMVPN [[DMVPN](#)] are used to establish tunnels among SD-WAN endpoints.

NHRP was originally intended for ATM address resolution, and as a result it misses many attributes which are necessary for dynamic end point CPE registration to controller, such as:

- Location identifier, such as Site Identifier, System ID, and/or Port ID.
- CPE attached GW information. When a CPE is instantiated within Cloud DC, the Cloud DC operator' GW to which the CPE is attached.
- Private <-> Public address mapping, which is needed when the CPEs use private addresses.
- IPsec configuration parameters (from controller to CPEs)

4. Gap Analysis in aggregating VPN paths and Internet paths

Most likely, enterprises, especially large ones, already have their CPEs interconnected by provider VPNs, such as EVPN, L2VPN, or L3VPN. The L2VPN or L3VPN can also be formed among all the CPEs directly attached to PEs, which is referred to as CPE based VPN as shown in the following diagram. The commonly used CPE based VPNs have CPE directly attached to PEs via VLANs (Ethernet). Therefore, the communication is secure. The BGP is used to distribute routes among CPEs.



=== or \\ indicates control plane communications

Figure 1: L2 or L3 VPNs over IP WAN

To use SD-WAN to aggregate Internet paths with the VPN paths, the CPEs need to have some ports connected to PEs and other Ports connected to the internet. NHRP & DSVPN/DMVPN can be used for the CPEs to be registered with their SD-WAN Controllers to establish secure tunnels among relevant CPEs.

That means the CPEs need to participate in two separate control planes: EVPN&BGP for CPE based VPN via links directly attached to PEs and NHRP & DSVPN/DMVPN. Two separate control planes not only add complexity to CPEs, but also increase operational cost.

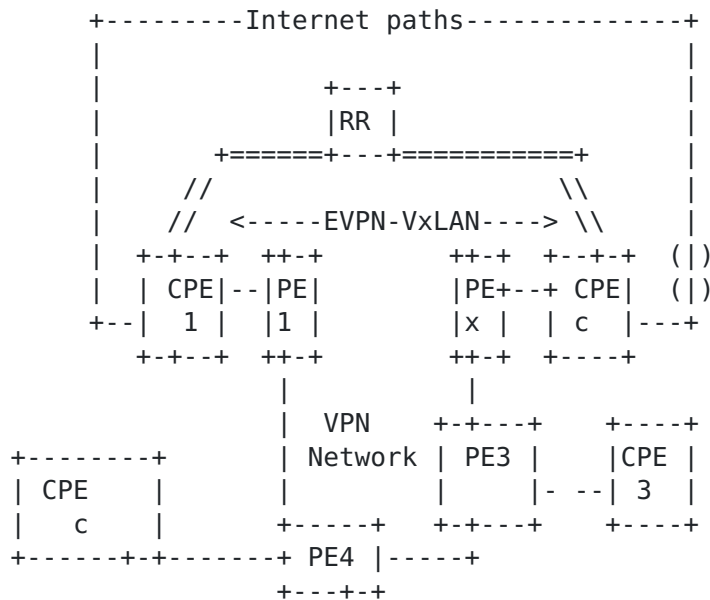


Figure 2: CPEs interconnected by VPN paths and Internet Paths

4.1. Gap analysis of Using BGP to cover SD-WAN paths

Since BGP is widely deployed, it is desirable to consider using BGP to control the SD-WAN paths instead of NHRP, DSVPN/DMVPN. This section analyzes the gaps of using BGP to control SD-WAN.

[RFC5512](#) and [Tunnel-Encap] describe methods for end points to advertise tunnel information and to trigger tunnel Establishment. [RFC5512](#) & [Tunnel-Encaps] have the Endpoint Address to indicate IPv4 or IPv6 address format Tunnel Encapsulation attribute to indicate different encapsulation formats, such as L2TPv3, GRE, VxLAN, IP in IP, etc. There are sub-TLVs to describe the detailed tunnel information for each of the encapsulations.

There is also Color sub-TLV to describe customer specified information about the tunnels (which can be creatively used for SD-

To express supporting multiple Encap types, multiple Extended communities with SAFI value = 7 can be used.

Here are some of the gaps using [RFC5512](#) and [Tunnel-Encap] to control SD-WAN:

- Doesn't have fields to carry detailed information of the remote CPE: such as Site-ID, System-ID, Port-ID

- Does not have the proper field to express IPsec configuration information from "Controller" (which can be RR) to CPEs.
- Does not have proper way for two peer CPEs to negotiate IPsec key based on the configuration sent from Controller.
- UDP NAT private address <-> public address mapping
- CPEs tend to communicate with a few other CPEs, not all the CPEs need to form mesh connections. Using BGP, CPEs can easily get dumped with too much information of other CPEs that they never need to communicate. NHRP only sends the relevant information for the interested end points for establishing tunnels. Therefore, need some form of "Registration" methods.

[VPN-over-Internet] describes a way to securely interconnect CPEs via IPsec using BGP. This method is useful, however, it still miss some aspects to aggregate CPE based VPN paths with internet paths that interconnect the CPEs. In addition:

- The draft assumes that CPE "register" with the RR. However, it does not say how. Should "NHRP" (modified version) be considered? In SD-WAN, Zero Touch Provisioning is expected. It is not acceptable to require manual configuration on RR which CPEs are controlled.

- The draft assumes that CPE and RR are connected by IPsec tunnel.

With zero touch provisioning, we need an automatic way to synchronize the IPsec SA between CPE and RR. The draft assumes:

A CPE must also be provisioned with whatever additional information is needed in order to set up an IPsec SA with each of the red RRs

- IPsec requires periodic refreshment of the keys. How to synchronize the refreshment among multiple nodes?

- IPsec usually only send configuration parameters to two end points

and let the two end points to negotiate the KEY. Now we assume that RR is responsible for creating the KEY for all end points. When one end point is confiscated, all other connections are impacted.

4.2. Gaps in preventing attacks to CPEs from their Internet ports

When CPEs have ports facing internet, it brings in the security risks of potential DDoS attacks to the CPEs from the ports facing internet. I.e. the CPE resource are attacked by unwanted traffic.

To mitigate security risk, it is absolutely necessary to enable Anti-DDoS feature on those CPEs to prevent major DDoS attack.

5. Gap analysis of CPEs not directly connected to VPN PEs

Because of the ephemeral property of the selected Cloud DCs, an enterprise or its network service provider may not have the direct links to the Cloud DCs that are optimal for hosting the enterprise's specific workloads/Apps. Under those circumstances, SD-WAN is a very flexible choice to interconnect the enterprise on-premises data centers & branch offices to its desired Cloud DCs.

However, SD-WAN paths over public internet can have unpredictable performance, especially over long distances and cross state/country boundaries. Therefore, it is highly desirable to place as much as possible the portion of SD-WAN paths over service provider VPN (e.g. enterprise's existing VPN) that have guaranteed SLA to minimize the distance/segments over public internet.

MEF Cloud Service Architecture [MEF-Cloud] also describes a use case of network operators needing to use SD-WAN over LTE or public internet for the last mile access that they do not have physical infrastructure.

Under those scenarios, one or both of the SD-WAN end points may not directly attached to the PEs of a SR Domain.

Using SD-WAN to connect the enterprise existing sites with the workloads in Cloud DC, the enterprise existing sites' CPEs have to be upgraded to support SD-WAN. If the workloads in Cloud DC need to be connected to many sites, the upgrade process can be very expensive.

[Net2Cloud-Problem] describes a hybrid network approach that integrates SD-WAN with traditional MPLS-based VPNs, to extend the existing MPLS-based VPNs to the Cloud DC Workloads over the access paths that are not under the VPN provider control. To make it working properly, a small number of the PEs of the MPLS VPN can be designated to connect to the remote workloads via SD-WAN secure IPsec tunnels. Those designated PEs are shown as fPE (floating PE or smart PE) in Figure below. Once the secure IPsec tunnels are established, the workloads in Cloud DC can be reached by the enterprise's VPN without upgrading all of the enterprise's existing

CPEs. The only CPE that needs to support SD-WAN would be a virtualized CPE instantiated within the cloud DC.

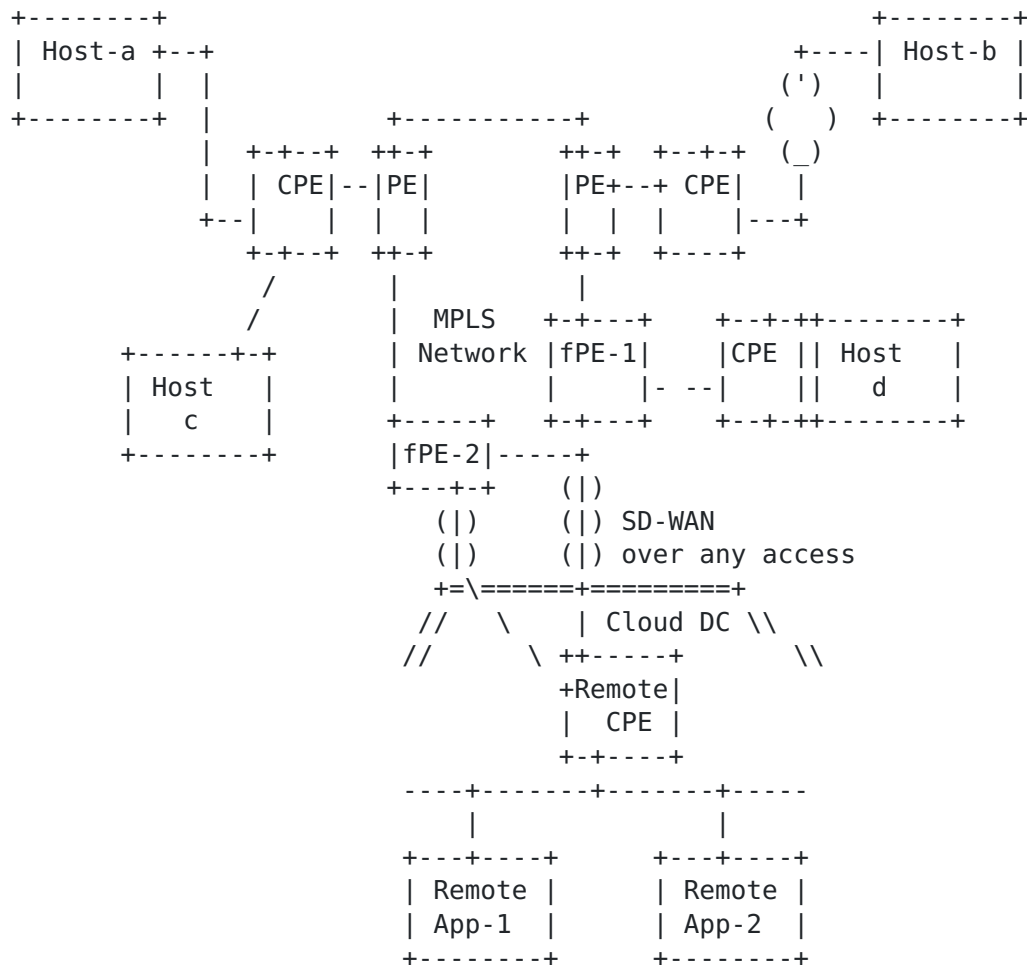


Figure 3: VPN Extension to Cloud DC

In Figure 3 above, the optimal Cloud DC to host the workloads (due to proximity, capacity, pricing, or other criteria chosen by the enterprises) does not happen to have a direct connection to the PEs of the MPLS VPN that interconnects the enterprise's existing sites.

5.1. Gap Analysis of Floating PEs to connect to Remote CPEs

To extend MPLS VPN to remote CPEs, it is necessary to establish secure tunnels (such as IPsec tunnels) between the Floating PEs and the remote CPEs.

Gap:

Even though a set of PEs can be manually selected to act as the floating PEs for a specific cloud data center, there are no standard protocols for those PEs to interact with the remote CPEs (most likely virtualized) instantiated in the third party cloud data centers (such as exchanging performance information or route information).

When there is more than one fPE available for use (as there should be for resiliency or the ability to support multiple cloud DCs scattered geographically), it is not straight to designate egress fPE to remote CPEs based on applications. There are too much applications traffic traversing PEs, it is not feasible for PEs to recognize applications carried by the payload.

5.2. NAT Traversing

Most cloud DCs only assign private IP addresses to the workloads instantiated. Therefore, the traffic to/from the workload usually need to traverse NAT.

5.3. Complication of use BGP between PE and remote CPEs via Internet

Even though EBGp (external BGP) Multihop method can be used to connect peers that are not directly connected to each other, there are still some complications/gaps in extending BGP from MPLS VPN PEs to remote CPEs via any access paths (e.g. internet):

EBGP Multi-hop scheme requires static configuration on both peers. To use EBGp between a PE and remote CPEs, the PE has to be statically configured with "next-hop" to the IP addresses of the CPEs. When remote CPEs, especially remote virtualized CPEs

dynamically instantiated or removed, the configuration on the PE Multi-Hop EBGP has to be changed accordingly.

Gap:

Egress peering engineering (EPE) is not enough. Running BGP on virtualized CPE in Cloud DC requires GRE tunnels being established first, which requires address and key management for the remote CPEs. [RFC 7024](#) (Virtual Hub & Spoke) and Hierarchical VPN is not enough

Also need a method to automatically trigger configuration changes on PE when remote CPEs' are instantiated or moved (IP address change) or deleted.

EBGP Multi-hop scheme does not have embedded security mechanism. The PE and remote CPEs needs secure communication channel when connected via public internet.

Remote CPEs, if instantiated in Cloud DC, might have to traverse NAT to reach PE. It is not clear how BGP can be used between devices outside the NAT and the entities behind the NAT. It is not clear how to configure the Next Hop on the PEs to reach private addresses.

[5.4. Designated Forwarder to the remote edges](#)

Among multiple floating PEs available for a remote CPE, multicast traffic from the remote CPE towards the MPLS VPN can be broadcasted back to the remote CPE due to the PE receiving the broadcast data frame forwarding the multicast/broadcast frame to other PEs that in turn send to all attached CPEs. This process may cause a traffic loop.

Therefore, it is necessary to designate one floating PE as the CPE's Designated Forwarder, similar to TRILL's Appointed Forwarders [[RFC6325](#)].

Gap: the MPLS VPN does not have features like TRILL's Appointed Forwarders.

5.5. Traffic Path Management

When there are multiple floating PEs that have established IPsec tunnels to the remote CPE, the remote CPE can forward the outbound traffic to the Designated Forwarder PE, which in turn forwards the traffic to egress PEs to the destinations. However, it is not straightforward for the egress PE to send back the return traffic to the Designated Forwarder PE.

Example of Return Path management using Figure 3 above.

- fPE-1 is desired for communication between App-1 <-> Host-a due to latency, pricing or other criteria.
- fPE-2 is desired for communication between App-1 <-> Host-b.

6. Manageability Considerations

TBD

7. Security Considerations

The intention of this draft is to identify the gaps in current and proposed SD-WAN approaches to the requirements identified in [Net2Cloud-problem].

Several of these approaches have gaps in meeting enterprise security requirements when tunneling their traffic over the Internet, as is the general intention of SD-WAN. See the individual sections above for further discussion of these security gaps.

8. IANA Considerations

This document requires no IANA actions. RFC Editor: Please remove this section before publication.

9. References

9.1. Normative References

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

9.2. 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.

[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

10. Acknowledgments

Acknowledgements to xxx for his review and contributions.

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

Authors' Addresses

Linda Dunbar
Huawei
Email: Linda.Dunbar@huawei.com

Andrew G. Malis
Huawei
Email: agmalis@gmail.com