

L3VPN WG
Internet Draft
Expiration Date: October 2004

Hamid Ould-Brahim
Nortel Networks

Eric C. Rosen
Cisco Systems

Yakov Rekhter
Juniper Networks

(Editors)

April 2004

**Using BGP as an Auto-Discovery
Mechanism for Layer-3 and Layer-2 VPNs**

[draft-ietf-l3vpn-bgpvpn-auto-02.txt](#)

Status of this Memo

This document is an Internet-Draft and is in full conformance with all provisions of [Section 10 of RFC2026](#) [[RFC-2026](#)].

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>.

Abstract

In any Provider Provisioned-Based VPN (PPVPN) scheme, the Provider Edge (PE) devices attached to a common VPN must exchange certain information as a prerequisite to establish VPN-specific connectivity. The purpose of this draft is to define a BGP based auto-discovery mechanism for both layer-2 VPN architectures and layer-3 VPNs ([[VPN-VR](#)]). This mechanism is based on the approach used by [[RFC2547-bis](#)] for distributing VPN routing information within the service provider(s). Each VPN scheme uses the mechanism

to automatically discover the information needed by that particular scheme.

1. Introduction

In any Provider Provisioned-Based VPN (PPVPN) scheme, the Provider Edge (PE) devices attached to a common VPN must exchange certain information as a prerequisite to establish VPN-specific connectivity. The purpose of this draft is to define a BGP based auto-discovery mechanism for both layer-2 VPN architectures (i.e., [[L2VPN-KOMP](#)], [[L2VPN-ROSEN](#)]) and layer-3 VPNs ([[VPN-VR](#)]). This mechanism is based on the approach used by [[RFC2547-bis](#)] for distributing VPN routing information within the service provider(s). Each VPN scheme uses the mechanism to automatically discover the information needed by that particular scheme.

In [[RFC2547-bis](#)] based layer-3 VPNs, VPN-specific routes are exchanged, along with the information needed to enable a PE to determine which routes belong to which VRFs. In [[VPN-VR](#)], virtual router (VR) addresses must be exchanged, along with the information needed to enable the PEs to determine which VRs are in the same VPN ("membership"), and which of those VRs are to have VPN connectivity ("topology"). Once the VRs are reachable through the tunnels, routes ("reachability") are then exchanged by running existing routing protocols per VPN basis.

The BGP-4 multiprotocol extensions are used to carry various information about VPNs for both layer-2 and layer-3 VPN architectures. VPN-specific information associated with the NLRI is encoded either as attributes of the NLRI, or as part of the NLRI itself, or both.

2. Provider-Provisioned VPN Reference Model

Both the layer-2 and layer-3 vpns architectures are using a network reference model as illustrated in figure 1.

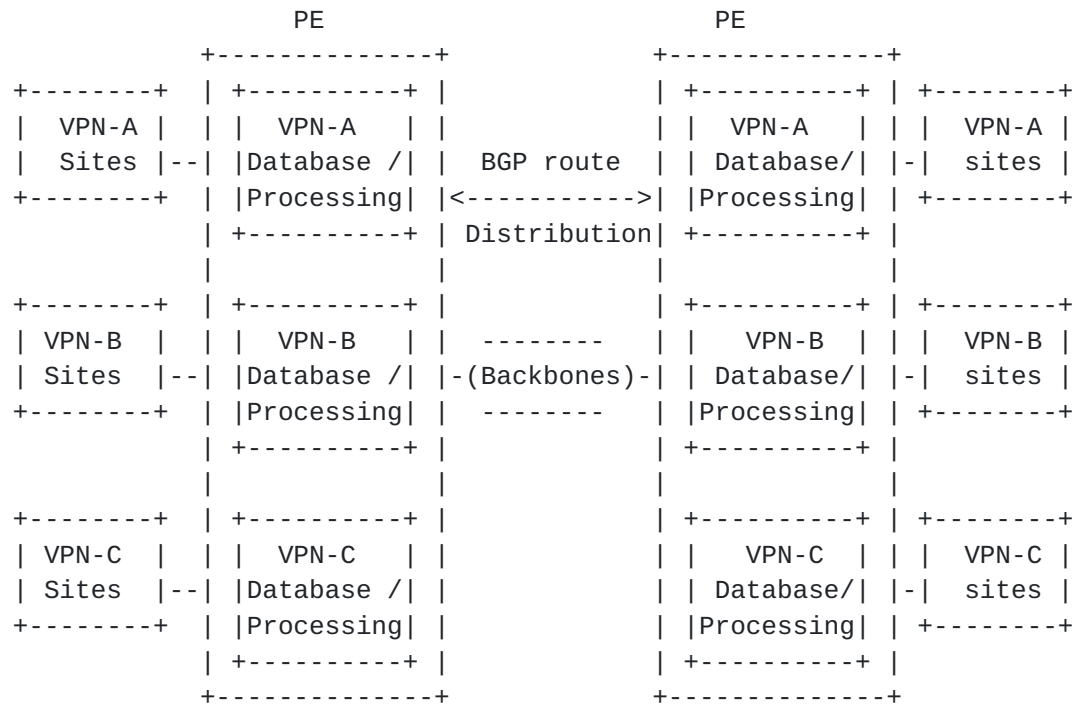


Figure 1: Network based VPN Reference Model

It is assumed that the PEs can use BGP to distribute information to each other. This may be via direct IBGP peering, via direct EBGP peering, via multihop BGP peering, through intermediaries such as Route Reflectors, through a chain of intermediate BGP connections, etc. It is assumed also that the PE knows what architecture it is supporting.

3. Carrying VPN information in BGP Multi-Protocol (BGP-MP) Attributes

The BGP-4 multiprotocol extensions are used to carry various information about VPNs for both layer-2 and layer-3 VPN architectures. VPN-specific information associated with the NLRI is encoded either as attributes of the NLRI, or as part of the NLRI itself, or both. The addressing information in the NLRI field is ALWAYS within the VPN address space, and therefore MUST be unique within the VPN. The address specified in the BGP next hop attribute, on the other hand, is in the service provider addressing space. In L3VPNs, the NLRI contains an address prefix which is within the VPN address space, and therefore must be unique within the VPN.

3.1 Carrying Layer-3 VPN Information in BGP-MP

This is done as follows. The NLRI is a VPN-IP address or a labeled

VPN-IP address.

In the case of the virtual router, the NLRI address prefix is an address of one of the virtual routers configured on the PE. Thus this mechanism allows the virtual routers to discover each other, to set up adjacencies and tunnels to each other, etc. In the case of [\[RFC2547-bis\]](#), the NLRI prefix represents a route to an arbitrary system or set of systems within the VPN.

[3.2](#) Carrying Layer-2 VPN Information in BGP-MP

The NLRI carries VPN layer-2 addressing information called VPN-L2 address. A VPN-L2 address is composed of a quantity beginning with an 8 bytes Route Distinguisher (RD) field and a variable length quantity encoded according to the layer-2 VPN architecture used.

Different layer-2 VPN solutions use the same common AFI, but different SAFI. The AFI indicates that the NLRI is carrying a VPN-l2 address, while the SAFI indicates solution-specific semantics and syntax of the VPN-l2 address that goes after the RD. The RD must be chosen so as it ensures that each NLRI is globally unique (i.e., the same NLRI does not appear in two VPNs).

BGP Route target extended community is used to constrain route distribution between PEs. The BGP Next hop carries the service provider tunnel endpoint address.

This draft doesn't preclude the use of additional extended communities for encoding specific l2vpn parameters.

[4](#). Interpretation of VPN Information in Layer-3 VPNs

[4.1](#) Interpretation of VPN Information in the [\[RFC2547-bis\]](#) Model

For details see [\[RFC2547-bis\]](#).

[4.2](#) Interpretation of VPN Information in the [\[VPN-VR\]](#) Model

[4.2.1](#) Membership Discovery

The VPN-ID format as defined in [\[RFC-2685\]](#) is used to identify a VPN. All virtual routers that are members of a specific VPN share the same VPN-ID. A VPN-ID is carried in the NLRI to make addresses of VRs globally unique. Making these addresses globally unique is necessary if one uses BGP for VRs' auto-discovery.

[4.2.1.1](#) Encoding of the VPN-ID in the NLRI

For the virtual router model, the VPN-ID is carried within the route distinguisher (RD) field. In order to hold the 7-bytes VPN-ID, the

first byte of RD type field is used to indicate the existence of the
Ould-Brahim & Rosen & Rekhter April 2004 [Page 4]

VPN-ID format. A value of 0x80 in the first byte of RD's type field indicates that the RD field is carrying the VPN-ID format. In this case, the type field range 0x8000-0x80ff will be reserved for the virtual router case.

4.2.1.2 VPN-ID Extended Community

A new extended community is used to carry the VPN-ID format. This attribute is transitive across the Autonomous system boundary. The type field of the VPN-ID extended community is of regular type to be assigned by IANA [[BGP-COMM](#)]. The remaining 7 bytes hold the VPN-ID value field as per [[RFC-2685](#)]. The BGP UPDATE message will carry information for a single VPN. It is the VPN-ID Extended Community, or more precisely route filtering based on the Extended Community that allows one VR to find out about other VRs in the same VPN.

4.2.2 VPN Topology Information

A new extended community is used to indicate different VPN topology values. This attribute is transitive across the Autonomous system boundary. The value of the type field for extended type is assigned by IANA. The first two bytes of the value field (of the remaining 6 bytes) are reserved. The actual topology values are carried within the remaining four bytes. The following topology values are defined:

Value	Topology Type
1	"Hub"
2	"Spoke"
3	"Mesh"

Arbitrary values can also be used to allow specific topologies to be constructed. VPN connectivity between two VRs within the same VPN is achieved if and only if at least one of them is a hub (the other is a hub or a spoke), or if both VRs are part of a full mesh VPN topology.

5. Interpretation of VPN Information in Layer-2 VPNs

The interpretation of the VPN information carried in the VPN-L2 address is to be specified as part of each L2VPN solution standardized by L2VPN working group.

6. Tunnel Discovery

Layer-3 VPNs and Layer-2 VPNs must be implemented through some form of tunneling mechanism, where the packet formats and/or the

addressing used within the VPN can be unrelated to that used to

route the tunneled packets across the backbone. There are numerous tunneling mechanisms that can be used by a network based VPN (e.g., IP/IP [[RFC-2003](#)], GRE tunnels [[RFC-1701](#)], IPSec [[RFC-2401](#)], and MPLS tunnels [[RFC-3031](#)]). Each of these tunnels allows for opaque transport of frames as packet payload across the backbone, with forwarding disjoint from the address fields of the encapsulated packets. A provider edge router may terminate multiple type of tunnels and forward packets between these tunnels and other network interfaces in different ways.

BGP can be used to carry tunnel endpoint addresses between edge routers. For scalability purposes, this draft recommends the use of tunneling mechanisms with demultiplexing capabilities such as IPSec, MPLS, and GRE (with respect to using GRE -the key field, it is no different than just MPLS over GRE, however there is no specification on how to exchange the key field, while there is a specification and implementations on how to exchange the label). Note that IP in IP doesn't have demultiplexing capabilities.

The BGP next hop will carry the service provider tunnel endpoint address. As an example, if IPSec is used as tunneling mechanism, the IPSec tunnel remote address will be discovered through BGP, and the actual tunnel establishment is achieved through IPSec signaling protocol.

When MPLS tunneling is used, the label carried in the NLRI field is associated with an address of a VR, where the address is carried in the NLRI and is encoded as a VPN-IP address.

7. Auto-Discovery and VR-[\[RFC2547-bis\]](#) Interworking Scenarios

Two interworking scenarios are considered when the network is using both virtual routers and [\[RFC2547-bis\]](#). The first scenario is a CE-PE relationship between a PE (implementing [\[RFC2547-bis\]](#)), and a VR appearing as a CE to the PE. The connection between the VR, and the PE can be either direct connectivity, or through a tunnel (e.g., IPSec).

The second scenario is when a PE is implementing both architectures. In this particular case, a single BGP session configured on the service provider network can be used to advertise either [\[RFC2547-bis\]](#) VPN information or the virtual router related VPN information. From the VR and the [\[RFC2547-bis\]](#) point of view there is complete separation from data path and addressing schemes. However the PE's interfaces are shared between both architectures.

A PE implementing only [\[RFC2547-bis\]](#) will not import routes from a BGP UPDATE message containing the VPN-ID extended community. On the other hand, a PE implementing the virtual router architecture will

not import routes from a BGP UPDATE message containing the route target extended community attribute.

The granularity at which the information is either [[RFC2547-bis](#)] related or VR-related is per BGP UPDATE message. Different SAFI numbers are used to indicate that the message carried in BGP multiprotocol extension attributes is to be handled by the VR or [[RFC2547-bis](#)] architectures. SAFI number of 128 is used for [[RFC2547-bis](#)] related format. A value of 129 for the SAFI number is for the virtual router (where the NLRI are carrying a labeled prefixes), and a SAFI value of 140 is for non labeled addresses.

8. Scalability Considerations

In this section, we briefly summarize the main characteristics of our model with respect to scalability.

Recall that the Service Provider network consists of (a) PE routers, (b) BGP Route Reflectors, (c) P routers (which are neither PE routers nor Route Reflectors), and, in the case of multi-provider VPNs, and (d) ASBRs.

A PE router, unless it is a Route Reflector should not retain VPN-related information unless it has at least one VPN with an Import Target identical to one of the VPN-related information Route Target attributes. Inbound filtering should be used to cause such information to be discarded. If a new Import Target is later added to one of the PE's VPNs (a "VPN Join" operation), it must then acquire the VPN-related information it may previously have discarded.

This can be done using the refresh mechanism described in [[BGP-RFSH](#)]. The outbound route filtering mechanism of [[BGP-ORE](#)] can also be used to advantage to make the filtering more dynamic.

Similarly, if a particular Import Target is no longer present in any of a PE's VPNs (as a result of one or more "VPN Prune" operations), the PE may discard all VPN-related information which, as a result, no longer have any of the PE's VPN's Import Targets as one of their Route Target Attributes.

Note that VPN Join and Prune operations are non-disruptive, and do not require any BGP connections to be brought down, as long as the refresh mechanism of [[BGP-RFSH](#)] is used.

As a result of these distribution rules, no one PE ever needs to maintain all routes for all VPNs; this is an important scalability consideration.

Route reflectors can be partitioned among VPNs so that each partition carries routes for only a subset of the VPNs supported by the Service Provider. Thus no single route reflector is required to maintain VPN-related information for all VPNs.

For inter-provider VPNs, if multi-hop EBGp is used, then the ASBRs need not maintain and distribute VPN-related information at all.

P routers do not maintain any VPN-related information. In order to properly forward VPN traffic, the P routers need only maintain routes to the PE routers and the ASBRs.

As a result, no single component within the Service Provider network has to maintain all the VPN-related information for all the VPNs. So the total capacity of the network to support increasing numbers of VPNs is not limited by the capacity of any individual component.

An important consideration to remember is that one may have any number of INDEPENDENT BGP systems carrying VPN-related information. This is unlike the case of the Internet, where the Internet BGP system must carry all the Internet routes. Thus one significant (but perhaps subtle) distinction between the use of BGP for the Internet routing and the use of BGP for distributing VPN-related information, as described in this document is that the former is not amenable to partition, while the latter is.

9. Security Considerations

This document describes a BGP-based auto-discovery mechanism which enables a PE router that attaches to a particular VPN to discover the set of other PE routers that attach to the same VPN. Each PE router that is attached to a given VPN uses BGP to advertise that fact. Other PE routers which attach to the same VPN receive these BGP advertisements. This allows that set of PE routers to discover each other. Note that a PE will not always receive these advertisements directly from the remote PEs; the advertisements may be received from "intermediate" BGP speakers.

It is of critical importance that a particular PE should not be "discovered" to be attached to a particular VPN unless that PE really is attached to that VPN, and indeed is properly authorized to be attached to that VPN. If any arbitrary node on the Internet could start sending these BGP advertisements, and if those advertisements were able to reach the PE routers, and if the PE routers accepted those advertisements, then anyone could add any site to any VPN. Thus the auto-discovery procedures described here presuppose that a particular PE trusts its BGP peers to be who they appear to be, and further that it can trust those peers to be properly securing their local attachments. (That is, a PE must trust that its peers are attached to, and are authorized to be attached to, the VPNs to which they claim to be attached.).

If a particular remote PE is a BGP peer of the local PE, then the

BGP authentication procedures of [RFC 2385](#) can be used to ensure that

the remote PE is who it claims to be, i.e., that it is a PE that is trusted.

If a particular remote PE is not a BGP peer of the local PE, then the information it is advertising is being distributed to the local PE through a chain of BGP speakers. The local PE must trust that its peers only accept information from peers that they trust in turn, and this trust relation must be transitive. BGP does not provide a way to determine that any particular piece of received information originated from a BGP speaker that was authorized to advertise that particular piece of information. Hence the procedures of this document should be used only in environments where adequate trust relationships exist among the BGP speakers.

Some of the VPN schemes which may use the procedures of this document can be made robust to failures of these trust relationships. That is, it may be possible to keep the VPNs secure even if the auto-discovery procedures are not secure. For example, a VPN based on the VR model can use IPsec tunnels for transmitting data and routing control packets between PE routers. An illegitimate PE router which is discovered via BGP will not have the shared secret which makes it possible to set up the IPsec tunnel, and so will not be able to join the VPN. Similarly, [[IPSEC-2547](#)] describes procedures for using IPsec tunnels to secure VPNs based on the [[RFC2547-bis](#)] model. The details for using IPsec to secure a particular sort of VPN depend on that sort of VPN and so are out of scope of the current document.

10. IANA Considerations

New AFI value to be assigned by IANA to indicate that the NLRI is carrying VPN-L2 Address as described in [section 3.2](#) to be used by all L2VPN solutions.

SAFI number of "128" is used for [[RFC2547-bis](#)].

SAFI number "129" for indicating that the NLRI is carrying information for VR-based solution.

SAFI number "140" for indicating that the NLRI is carrying information for VR for non labeled prefixes.

New Extended Community to be assigned by IANA and used for Topology values for VR-based L3VPN solution see [section 4.2.2](#).

New Extended Community to be assigned by IANA for carrying VPN-ID format based on [RFC2685](#) format (see [section 4.2.1.2](#))

11. Use of BGP Capability Advertisement

A BGP speaker that uses VPN information as described in this document with multiprotocol extensions should use the Capability

Advertisement procedures [[RFC-3392](#)] to determine whether the speaker could use Multiprotocol Extensions with a particular peer.

[12. Normative References](#)

[BGP-COMM] Ramachandra, Tappan, et al., "BGP Extended Communities Attribute", June 2001, work in progress

[BGP-MP] Bates, Chandra, Katz, and Rekhter, "Multiprotocol Extensions for BGP4", February 1998, [RFC 2283](#)

[RFC-3107] Rekhter Y, Rosen E., "Carrying Label Information in BGP4", January 2000, [RFC3107](#)

[RFC2547-bis] Rosen E., et al, "BGP/MPLS VPNs", Work in Progress.

[RFC-2685] Fox B., et al, "Virtual Private Networks Identifier", [RFC 2685](#), September 1999.

[RFC-3392] Chandra, R., et al., "Capabilities Advertisement with BGP-4", [RFC3392](#), May 2002.

[VPN-VR] Knight, P., Ould-Brahim H., Gleeson, B., "Network based IP VPN Architecture using Virtual Routers", Work in Progress.

[13. Informative References](#)

[L2VPN-ROSEN] Rosen, E., Radoaca, V., "Provisioning Models and Endpoint Identifiers in L2VPN Signaling", Work in Progress.

[L2VPN-KOMP] Kompella, K., et al., "Virtual Private LAN Service", Work in Progress.

[L2VPN-VKOMP-LASS] Kompella, V., Lasserre, M., et al., "Transparent VLAN Services over MPLS", Work in Progress.

[RFC-1701] Hanks, S., Li, T., Farinacci, D. and P. Traina, "Generic Routing Encapsulation (GRE)", [RFC 1701](#), October 1994.

[RFC-2003] Perkins, C., "IP Encapsulation within IP", [RFC2003](#), October 1996.

[RFC-2026] Bradner, S., "The Internet Standards Process -- Revision 3", [RFC2026](#), October 1996.

[RFC-2401] Kent S., Atkinson R., "Security Architecture for the Internet Protocol", [RFC2401](#), November 1998.

[RFC-2119] Bradner, S., "Key words for use in RFCs to Indicate

Requirement Levels", [RFC 2119](#), March 1997.

Ould-Brahim & Rosen & Rekhter

April 2004

[Page 10]

[TLS-TISSA] "BGP/MPLS Layer-2 VPN", [draft-tsenevir-bgp12vpn-01.txt](#), work in progress, July 2001.

[IPSEC-2547] Rosen, E., et al., "Use of PE-PE IPsec in [RFC2547](#) VPNs", Work in Progress.

[BGP-RFSH] Chen, A., "Route Refresh Capability for BGP-4", [RFC2918](#), September 2000.

[BGP-ORF] Chen, E., and Rekhter, Y., "Cooperative Route Filtering Capability for BGP-4", Work in Progress.

14. Intellectual Property Rights Notices

The IETF takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Information on the IETF's procedures with respect to rights in standards-track and standards-related documentation can be found in [BCP-11](#). Copies of claims of rights made available for publication and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementors or users of this specification can be obtained from the IETF Secretariat.

15. Contributors

Bryan Gleeson
Tahoe Networks
3052 Orchard Drive
San Jose, CA 95134 USA
Email: bryan@tahoenetworks.com

Peter Ashwood-Smith
Nortel Networks
P.O. Box 3511 Station C,
Ottawa, ON K1Y 4H7, Canada
Phone: +1 613 763 4534
Email: petera@nortelnetworks.com

Luyuan Fang
AT&T
200 Laurel Avenue
Middletown, NJ 07748
Email: Luyuanfang@att.com

Phone: +1 (732) 420 1920

Jeremy De Clercq
Alcatel
Francis Wellesplein 1
B-2018 Antwerpen, Belgium
Phone: +32 3 240 47 52
Email: jeremy.de_clercq@alcatel.be

Riad Hartani
Caspian Networks
170 Baytech Drive
San Jose, CA 95143
Phone: 408 382 5216
Email: riad@caspiannetworks.com

Tissa Senevirathne
Force10 Networks
1440 McCarthy Blvd,
Milpitas, CA 95035.

Phone: 408-965-5103
Email: tsenevir@hotmail.com

16. Authors Information

Hamid Ould-Brahim
Nortel Networks
P O Box 3511 Station C
Ottawa, ON K1Y 4H7, Canada
Email: hbrahim@nortelnetworks.com

Eric C. Rosen
Cisco Systems, Inc.
1414 Massachusetts Avenue
Boxborough, MA 01719
E-mail: erosen@cisco.com

Yakov Rekhter
Juniper Networks
1194 N. Mathilda Avenue
Sunnyvale, CA 94089
Email: yakov@juniper.net

Full Copyright Statement

Copyright (C) The Internet Society (2004). All Rights Reserved. This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

