

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: August 19, 2014

Yiqun Cai
Microsoft
Sri Vallepalli
Heidi Ou
Cisco Systems, Inc.
Andy Green
British Telecom
February 15, 2014

PIM Designated Router Load Balancing
draft-ietf-pim-drlb-03.txt

Abstract

On a multi-access network, one of the PIM routers is elected as a Designated Router (DR). On the last hop network, the PIM DR is responsible for tracking local multicast listeners and forwarding traffic to these listeners if the group is operated in PIM SM. In this document, we propose a modification to the PIM SM protocol that allows more than one of these last hop routers to be selected so that the forwarding load can be distributed to and handled among these routers.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on August 19, 2014.

Copyright Notice

Copyright (c) 2014 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. Terminology	3
2. Introduction	3
3. Applicability	6
4. Functional Overview	6
4.1. GDR Candidates	7
4.2. Hash Mask	7
4.3. PIM Hello Options	8
5. Hello Option Formats	9
5.1. PIM DR Load Balancing Capability (DRLBC) Hello Option	9
5.2. PIM DR Load Balancing GDR (DRLBGDR) Hello Option	10
6. Protocol Specification	11
6.1. PIM DR Operation	11
6.2. PIM GDR Candidate Operation	11
6.3. PIM Assert Modification	12
7. IANA Considerations	14
8. Security Considerations	14
9. Acknowledgement	14
10. References	14
10.1. Normative Reference	14
10.2. Informative References	14
Authors' Addresses	15

1. Terminology

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 [RFC2119].

With respect to PIM, this document follows the terminology that has been defined in [RFC4601].

This document also introduces the following new acronyms:

- o GDR: GDR stands for "Group Designated Router". For each multicast group, a hash algorithm (described below) is used to select one of the routers as a GDR. The GDR is responsible for initiating the forwarding tree building for the corresponding group.
- o GDR Candidate: a last hop router that has potential to become a GDR. A GDR Candidate must have the same DR priority and must run the same GDR election hash algorithm as the DR router. It must send and process received new PIM Hello Options as defined in this document. There might be more than one GDR Candidate on a LAN. But only one can become GDR for a specific multicast group.

2. Introduction

On a multi-access network such as an Ethernet, one of the PIM routers is elected as a DR. The PIM DR has two roles in the PIM protocol. On the first hop network, the PIM DR is responsible for registering an active source with the Rendezvous Point (RP) if the group is operated in PIM SM. On the last hop network, the PIM DR is responsible for tracking local multicast listeners and forwarding to these listeners if the group is operated in PIM SM.

Consider the following last hop network in Figure 1:

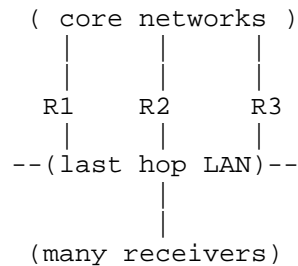


Figure 1: Last Hop Network

Assume R1 is elected as the Designated Router. According to [RFC4601], R1 will be responsible for forwarding to the last hop LAN. In addition to keeping track of IGMP and MLD membership reports, R1 is also responsible for initiating the creation of source and/or shared trees towards the senders or the RPs.

Forcing sole data plane forwarding responsibility on the PIM DR proves a limitation in the protocol. In comparison, even though an OSPF DR, or an IS-IS DIS, handles additional duties while running the OSPF or IS-IS protocols, they are not required to be solely responsible for forwarding packets for the network. On the other hand, on a last hop LAN, only the PIM DR is asked to forward packets while the other routers handle only control traffic (and perhaps drop packets due to RPF failures). The forwarding load of a last hop LAN is concentrated on a single router.

This leads to several issues. One of the issues is that the aggregated bandwidth will be limited to what R1 can handle towards this particular interface. These days, it is very common that the last hop LAN usually consists of switches that run IGMP/MLD or PIM snooping. This allows the forwarding of multicast packets to be restricted only to segments leading to receivers who have indicated their interest in multicast groups using either IGMP or MLD. The emergence of the switched Ethernet allows the aggregated bandwidth to exceed, some times by a large number, that of a single link. For example, let us modify Figure 1 and introduce an Ethernet switch in Figure 2.

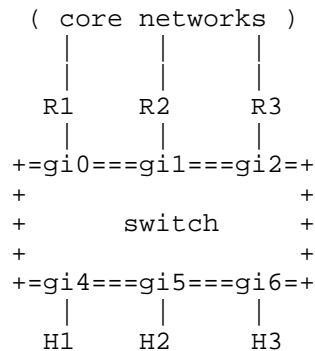


Figure 2: Last Hop Network with Ethernet Switch

Let us assume that each individual link is a Gigabit Ethernet. Each router, R1, R2 and R3, and the switch have enough forwarding capacity to handle hundreds of Gigabits of data.

Let us further assume that each of the hosts requests 500 mbps of data and different traffic is requested by each host. This represents a total 1.5 gbps of data, which is under what each switch or the combined uplink bandwidth across the routers can handle, even under failure of a single router.

On the other hand, the link between R1 and switch, via port gi0, can only handle a throughput of 1gbps. And if R1 is the only router, the PIM DR elected using the procedure defined by RFC 4601, at least 500 mbps worth of data will be lost because the only link that can be used to draw the traffic from the routers to the switch is via gi0. In other words, the entire network's throughput is limited by the single connection between the PIM DR and the switch (or the last hop LAN as in Figure 1).

The problem may also manifest itself in a different way. For example, R1 happens to forward 500 mbps worth of unicast data to H1, and at the same time, H2 and H3 each requests 300 mbps of different multicast data. Once again packet drop happens on R1 while in the mean time, there is sufficient forwarding capacity left on R2 and R3 and link capacity between the switch and R2/R3.

Another important issue is related to failover. If R1 is the only forwarder on the last hop network, in the event of a failure when R1 goes out of service, multicast forwarding for the entire network has to be rebuilt by the newly elected PIM DR. However, if there was a way that allowed multiple routers to forward to the network for different groups, failure of one of the routers would only lead to

disruption to a subset of the flows, therefore improving the overall resilience of the network.

In this document, we propose a modification to the PIM protocol that allows more than one of these routers, called Group Designated Router (GDR) to be selected so that the forwarding load can be distributed to and handled by a number of routers.

3. Applicability

The proposed change described in this specification applies to PIM SM last hop routers only.

It does not alter the behavior of a PIM DR on the first hop network. This is because the source tree is built using the IP address of the sender, not the IP address of the PIM DR that sends the registers towards the RP. The load balancing between first hop routers can be achieved naturally if an IGP provides equal cost multiple paths (which it usually does in practice). And distributing the load to do registering does not justify the additional complexity required to support it.

4. Functional Overview

In the existing PIM DR election, when multiple last hop routers are connected to a multi-access network (for example, an Ethernet), one of them is selected to act as PIM DR. The PIM DR is responsible for sending Join/Prune messages towards the RP or source. To elect the PIM DR, each PIM router on the network examines the received PIM Hello messages and compares its DR priority and IP address with those of its neighbors. The router with the highest DR priority is the PIM DR. If there are multiple such routers, their IP addresses are used as the tie-breaker, as described in [RFC4601].

In order to share forwarding load among last hop routers, besides the normal PIM DR election, the GDR is also elected on the last hop multi-access network. There is only one PIM DR on the multi-access network, but there might be multiple GDR Candidates.

For each multicast group, a hash algorithm is used to select one of the routers to be the GDR. Hash Masks are defined for Source, Group and RP separately, in order to handle PIM ASM/SSM. The masks are announced in PIM Hello by DR as a DR Load Balancing GDR (DRLBGDR) Hello Option. Besides that, a DR Load Balancing Capability (DRLBC) Hello Option, which contains hash algorithm type, is also announced by router interfaces which have this specification supported. Last

hop routers who are with the new DRLBC Option, and with the same GDR election hash algorithm and the same DR priority as the PIM DR are GDR Candidates.

A hash algorithm based on the announced Source, Group or RP masks allows one GDR to be assigned to a corresponding multicast group, and that GDR is responsible for initiating the creation of the multicast forwarding tree for the group.

4.1. GDR Candidates

GDR is the new concept introduced by this specification. GDR Candidates are routers eligible for GDR election on the LAN. To become a GDR Candidate, a router MUST support this specification, have the same DR priority and run the same GDR election hash algorithm as the DR on the LAN.

For example, assume there are 4 routers on the LAN: R1, R2, R3 and R4, which all support this specification on the LAN. R1, R2 and R3 have the same DR priority while R4's DR priority is less preferred. In this example, R4 will not be eligible for GDR election, because R4 will not become a PIM DR unless all of R1, R2 and R3 go out of service.

Further assume router R1 wins the PIM DR election, and R1, R2 run the same hash algorithm for GDR election, while R3 runs a different one. Then only R1 and R2 will be eligible for GDR election, R3 will not.

As a DR, R1 will include its own Load Balancing Hash Masks, and also the identity of R1 and R2 (the GDR Candidates) in its DRLBGDR Hello Option.

4.2. Hash Mask

A Hash Mask is used to extract a number of bits from the corresponding IP address field (32 for v4, 128 for v6), and calculate a hash value. A hash value is used to select a GDR from GDR Candidates advertised by PIM DR. For example, 0.255.0.0 defines a Hash Mask for an IPv4 address that masks the first, the third and the fourth octets.

There are three Hash Masks defined,

- o RP Hash Mask
- o Source Hash Mask
- o Group Hash Mask

The Hash Masks MUST be configured on the PIM routers that can

potentially become a PIM DR.

A simple Modulo hash algorithm will be discussed in this document. However, to allow other hash algorithm to be used, a 4-bytes "Hash Algorithm Type" field is included in DRLBC Hello Option to specify the hash algorithm used by a last hop router.

If different hash algorithm types are advertised among last hop routers, only last hop routers running the same hash algorithm as the DR (and having the same DR priority as the DR) are eligible for GDR election.

For ASM groups, a hash value is calculated using the following Modulo algorithm:

```
o  hashvalue_RP = (((RP_address & RP_hashmask) >> N) & 0xFFFF) % M
```

RP_address is the address of the RP defined for the group. N is the number of zeros, counted from the least significant bit of the RP_hashmask. For example, for a given IPv4 RP_hashmask 0.255.0.0, N will be 16. M is the number of GDR Candidates as described above.

If RP_hashmask is 0, a hash value is also calculated using the group Hash Mask in a similar fashion.

```
o  hashvalue_Group = (((Group_address & Group_hashmask) >> N) &
    0xFFFF) % M
```

For SSM groups, a hash value is calculated using both the source and group Hash Mask

```
o  hashvalue_SG = (((((Source_address & Source_hashmask) >> N_S) &
    0xFFFF) ^ (((Group_address & Group_hashmask) >> N_G) & 0xFFFF)) %
    M
```

4.3. PIM Hello Options

When a last hop PIM router sends a PIM Hello from an interface with this specification support, it includes a new option, called "Load Balancing Capability (DRLBC)".

Besides this DRLBC Hello Option, the elected PIM DR also includes a new "DR Load Balancing GDR (DRLBGDR) Hello Option". The DRLBGDR Hello Option consists of three Hash Masks as defined above and also the sorted addresses of all GDR Candidates on the last hop network.

The elected PIM DR uses DRLBC Hello Option advertised by all routers on the last hop network to compose its DRLBGDR . The GDR Candidates

use DRLBGDR Hello Option advertised by PIM DR to calculate hash value.

5. Hello Option Formats

5.1. PIM DR Load Balancing Capability (DRLBC) Hello Option

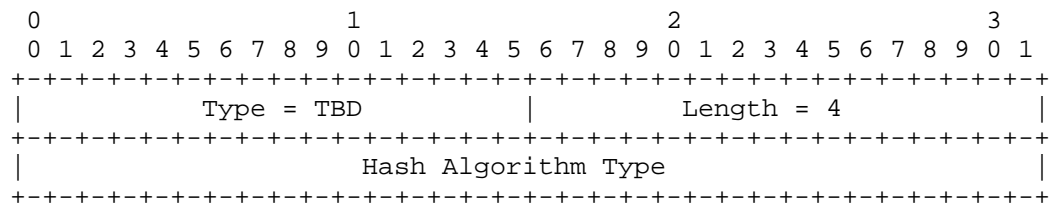


Figure 3: Capability Hello Option

Type: TBD.

Length: 4 octets

Hash Algorithm Type: 0 for Modulo hash algorithm

This DRLBC Hello Option SHOULD be advertised by last hop routers from interfaces which support this specification.

5.2. PIM DR Load Balancing GDR (DRLBGDR) Hello Option

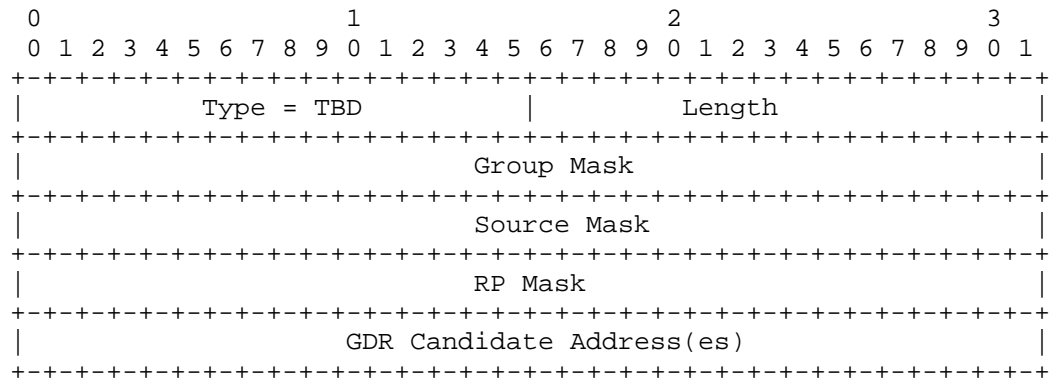


Figure 4: GDR Hello Option

Type: TBD

Length:

Group Mask (32/128 bits): Mask

Source Mask (32/128 bits): Mask

RP Mask (32/128 bits): Mask

All masks MUST be in the same address family, with the same length.

GDR Address (32/128 bits): Address(es) of GDR Candidate(s)

All addresses must be in the same address family. The addresses are sorted from high to low. The order is converted to the ordinal number associated with each GDR candidate in hash value calculation. For example, addresses advertised are R3, R2, R1, the ordinal number assigned to R3 is 0, to R2 is 1 and to R1 is 2. If "Interface ID" option (type 31) presents in a GDR Candidate's PIM Hello message, and the "Router ID" portion is non-zero,

* For IPv4, the "GDR Candidate Address" will be set directly to "Router ID".

* For IPv6, the "GDR Candidate Address" will be set to the IPv4-IPv6 translated address of "Router ID", as described in [RFC4291], that is the "Router-ID" is appended to the prefix of 96-bits zeros.

If the "Interface ID" option is not present in a GDR Candidate's PIM Hello message, or if the "Interface ID" option is present, but "Router ID" field is zero, the "GDR Candidate Address" will be the IPv4 or IPv6 source address from PIM Hello message.

This DRLBGDR Hello Option SHOULD only be advertised by the elected PIM DR.

6. Protocol Specification

6.1. PIM DR Operation

The DR election process is still the same as defined in [RFC4601]. A DR that has this specification enabled on the interface, advertises the new LBGRD Hello Option, which contains value of masks from user configuration, followed by a sorted list of addresses of all GDR Candidates. Moreover, same as non-DR routers, DR also advertises DRLBC Hello Option to indicate its capability of supporting this specification and the type of its GDR election hash algorithm.

If a PIM DR receives a neighbor Hello with DRLBGRD Option, the PIM DR SHOULD ignore the TLV.

If a PIM DR receives a neighbor DRLBC Hello Option, which contains the same hash algorithm type as the DR, and the neighbor has the same DR priority as the DR, PIM DR SHOULD consider the neighbor as a GDR Candidate and insert the neighbor's address into the sorted list of DRLBGRD Option.

6.2. PIM GDR Candidate Operation

When an IGMP join is received, without this proposal, router R1 (the PIM DR) will handle the join and potentially run into the issues described earlier. Using this proposal, a hash algorithm is used to determine which router is going to be responsible for building forwarding trees on behalf of the host.

The algorithm works as follows, assuming the router in question is X, which is a GDR Candidate, and its ordinal number assigned implicitly by PIM DR in DRLBGDR Hello Option is Ox:

- o If the group is ASM, and the RP Hash Mask announced by the PIM DR is not zero, calculate the value of hashvalue_RP. If hashvalue_RP is equal to Ox, X becomes the GDR.

For example, X with IPv4 address 10.1.1.3, receives a DRLBGDR Hello Option from the DR, which announces RP Hash Mask 0.255.0.0, and a

list of GDR Candidates, sorted by IP addresses from high to low, 10.1.1.3, 10.1.1.2 and 10.1.1.1. The ordinal number assigned to those addresses would be 0 for 10.1.1.3 (X), 1 for 10.1.1.2, and 2 for 10.1.1.1. Assume there are 2 RPs: RP1 172.3.10.10 for Group1 and RP2 172.2.10.10 for Group2. Following the modulo hash algorithm

$$\text{hashvalue_RP} = (((\text{RP_address} \& \text{RP_hashmask}) \gg N) \& 0\text{xFFFF}) \% M$$

Here N is 16 for 0.255.0.0, and M is 3 for the total number of GDR Candidates. The hashvalue_RP for RP1 172.3.10.10 is 0, matches the ordinal number assigned to X. X will be the GDR for Group1, which uses 172.3.10.10 as the RP. The hashvalue_RP for RP2 172.2.10.10 is 2, which is different from X's ordinal number, hence, X will not be GDR for Group2.

- o If the group is ASM, and the RP Hash Mask announced by the PIM DR is zero, obtain the value of hashvalue_Group. Compare hashvalue_Group with 0x, to decide if X is the GDR.
- o If the group is SSM, then use hashvalue_SG to determine if X is the GDR.

If X is the GDR for the group, X will be responsible for building the forwarding tree.

A router interface where this protocol is enabled advertises DRLBC Hello Option in its PIM Hello, even if the router may not be a GDR Candidate.

A GDR Candidate may receive a DRLBGDR Hello Option from PIM DR, with different Hash Masks from those configured on it, The GDR Candidate must use the Hash Masks advertised by the PIM DR to calculate the hash value.

A GDR Candidate may receive a DRLBGDR Hello Option from a non-DR PIM router. The GDR Candidate must ignore such DRLBGDR Hello Option.

A GDR Candidate may receive a Hello from the elected PIM DR, and the PIM DR does not support this specification. The GDR election described by this specification will not take place, that is only the PIM DR joins the multicast tree.

6.3. PIM Assert Modification

It is possible that the identity of the GDR might change in the middle of an active flow. Examples this could happen include:

1. When a new PIM router comes up

2. When a GDR restarts

When the GDR changes, existing traffic might be disrupted. Duplicates or packet loss might be observed. To illustrate the case, consider the following scenario: there are two streams G1 and G2. R1 is the GDR for G1, and R2 is the GDR for G2. When R3 comes up online, it is possible that R3 becomes GDR for both G1 and G2, hence R2 starts to build the forwarding tree for G1 and G2. If R1 and R2 stop forwarding before R3 completes the process, packet loss might occur. On the other hand, if R1 and R2 continue forwarding while R3 is building the forwarding trees, duplicates might occur.

This is not a typical deployment scenario but it still might happen. Here we describe a mechanism to minimize the impact. The motivation is that we want to minimize packet loss. And therefore, we would allow a small amount of duplicates and depend on PIM Assert to minimize the duplication.

When the role of GDR changes as above, instead of immediately stopping forwarding, R1 and R2 continue forwarding to G1 and G2 respectively, while at the same time, R3 build forwarding trees for G1 and G2. This will lead to PIM Asserts.

Due to the introduction of GDR, this document suggests the following modification to the Assert packet: if a router enables this specification on its downstream interface, but it is not a GDR, it would adjust its Assert metric to (PIM_ASSERT_INFINITY - 1).

Using the above example, assume R1 and R3 agree on the new GDR, which is R3. R1 will set its Assert metric as (PIM_ASSERT_INFINITY - 1). That will make R3, which has normal metric in its Assert as the Assert winner.

For G2, assume it takes a little bit longer time for R2 to find out that R3 is the new GDR and still thinks itself being the GDR while R3 already has assumed the role of GDR. Since both R2 and R3 think they are GDRs, they further compare the metric and IP address. If R3 has the better routing metric, or same metric but better tie-breaker, the result will be consistent with GDR selection. If unfortunately, R2 has the better metric or same metric but better tie-breaker R2 will become the Assert winner and continues to forward traffic. This will continue until:

1. The next PIM Hello option from DR is seen that selects R3 as the GDR.
 2. R3 will build the forwarding tree and send an Assert.
- The process continues until R2 agrees to the selection of R3 as being the GDR, and set its own Assert metric to (PIM_ASSERT_INFINITY - 1), which will make R3 the Assert winner. During the process, we will see intermittent duplication of traffic but packet loss will be

minimized. In the unlikely case that R2 never relinquishes its role as GDR (while every other router thinks otherwise), the proposed mechanism also helps to keep the duplication to a minimum until manual intervention takes place to remedy the situation.

7. IANA Considerations

Two new PIM Hello Option Types are required to be assigned to the DR Load Balancing messages. [HELLO-OPT], this document recommends 34(0x22) as the new "PIM DR Load Balancing Capability Hello Option", and 35(0x23) as the new "PIM DR Load Balancing GDR Hello Option".

8. Security Considerations

Security of the PIM DR Load Balancing Hello message is only guaranteed by the security of PIM Hello message, so the security considerations for PIM Hello messages as described in PIM-SM [RFC4601] apply here.

9. Acknowledgement

The authors would like to thank Steve Simlo, Taki Millonis for helping with the original idea, Bill Atwood for review comments, Stig Venaas, Toerless Eckert and Rishabh Parekh for helpful conversation on the document.

10. References

10.1. Normative Reference

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4601] Fenner, B., Handley, M., Holbrook, H., and I. Kouvelas, "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)", RFC 4601, August 2006.

10.2. Informative References

- [RFC3973] Adams, A., Nicholas, J., and W. Siadak, "Protocol Independent Multicast - Dense Mode (PIM-DM): Protocol Specification (Revised)", RFC 3973, January 2005.
- [RFC5015] Handley, M., Kouvelas, I., Speakman, T., and L. Vicisano,

"Bidirectional Protocol Independent Multicast (BIDIR-PIM)", RFC 5015, October 2007.

[RFC6395] Gulrajani, S. and S. Venaas, "An Interface Identifier (ID) Hello Option for PIM", RFC 6395, October 2011.

[RFC4291] Hinden, R. and L. S., "IP Version 6 Addressing Architecture", RFC 6890, February 2006.

[HELLO-OPT]
IANA, "PIM Hello Options", PIM-HELLO-OPTIONS per RFC4601 <http://www.iana.org/assignments/pim-hello-options>, March 2007.

Authors' Addresses

Yiqun Cai
Microsoft
La Avenida
Mountain View, CA 94043
USA

Email: yiqunc@microsoft.com

Sri Vallepalli
Cisco Systems, Inc.
Tasman Drive
San Jose, CA 95134
USA

Email: svallepa@cisco.com

Heidi Ou
Cisco Systems, Inc.
Tasman Drive
San Jose, CA 95134
USA

Email: hou@cisco.com

Andy Green
British Telecom
Adastral Park
Ipswich IP5 2RE
United Kingdom

Email: andy.da.green@bt.com

