              Routing-Related Design Choices for IPv6 Networks
                   draft-ietf-v6ops-design-choices-12

Abstract

   This document presents advice on certain routing-related design
   choices that arise when designing IPv6 networks (both dual-stack and
   IPv6-only).  The intended audience is someone designing an IPv6
   network who is knowledgeable about best current practices around IPv4
   network design, and wishes to learn the corresponding practices for
   IPv6.

Status of This Memo

Copyright Notice

Table of Contents

1.  Introduction

   This document discusses routing-related design choices that arise
   when designing an IPv6-only or dual-stack network.  The focus is on
   choices that do not come up when designing an IPv4-only network.  The
   document presents each choice and the alternatives, and then
   discusses the pros and cons of the alternatives in detail.  Where
   consensus currently exists around the best practice, this is
   documented; otherwise the document simply summarizes the current
   state of the discussion.  Thus this document serves to both document
   the reasoning behind best current practices for IPv6, and to allow a
   designer to make an informed choice where no such consensus exists.

   The design choices presented apply to both Service Provider and
   Enterprise network environments.  Where choices have selection
   criteria which differ between the Service Provider and the Enterprise

environment, this is noted.  The designer is encouraged to ensure
that they familiarize themselves with any of the discussed
technologies to ensure the best selection is made for their
environment.

This document does not present advice on strategies for adding IPv6
to a network, nor does it discuss transition in these areas, see
[RFC6180] for general advice,[RFC6782] for wireline service
providers, [RFC6342]for mobile network providers, [RFC5963] for
exchange point operators, [RFC6883] for content providers, and both
[RFC4852] and [RFC7381] for enterprises.  Nor does this document
discuss the particulars of creating an IPv6 addressing plan; for
advice in this area, see [RFC5375] or [v6-addressing-plan].  The
document focuses on unicast routing design only and does not cover
multicast or the issues involved in running MPLS over IPv6 transport

Section 2 presents and discusses a number of design choices.
Section 3 discusses some general themes that run through these
choices.

## 2.  Design Choices

Each subsection below presents a design choice and discusses the pros
and cons of the various options.  If there is consensus in the
industry for a particular option, then the consensus position is
noted.

## 2.1.  Addresses

This section discusses the choice of addresses for router loopbacks
and links between routers.  It does not cover the choice of addresses
for end hosts.

In IPv6, an interface is always assigned a Link-Local Address (LLA)
[RFC4291].  The link-local address can only be used for communicating
with devices that are on-link, so often one or more additional
addresses are assigned which are able to communicate off-link.  This
additional address or addresses can be one of three types:

o  Provider-Independent Global Unicast Address (PI GUA): IPv6 address
   allocated by a regional address registry [RFC4291]

o  Provider-Aggregatable Global Unicast Address (PA GUA): IPv6
   Address allocated by your upstream service provider

o  Unique Local Address (ULA): IPv6 address locally assigned
   [RFC4193]

This document uses the term "multi-hop address" to collectively refer
to these three types of addresses.

PI GUAs are, for many situations, the most flexible of these choices.
Their main disadvantages are that a regional address registry will
only allocate them to organizations that meet certain qualifications,
and one must pay an annual fee.  These disadvantages mean that many
smaller organization may not qualify or be willing to pay for these
addresses.

PA GUAs have the advantage that they are usually provided at no extra
charge when you contract with an upstream provider.  However, they
have the disadvantage that, when switching upstream providers, one
must give back the old addresses and get new addresses from the new
provider ("renumbering").  Though IPv6 has mechanisms to make
renumbering easier than IPv4, these techniques are not generally
applicable to routers and renumbering is still fairly hard [RFC5887]
[RFC6879] [RFC7010] .  PA GUAs also have the disadvantage that it is
not easy to have multiple upstream providers ("multi-homing") if they
are used (see "Ingress Filtering Problem" in [RFC5220] ).

ULAs have the advantage that they are extremely easy to obtain and
cost nothing.  However, they have the disadvantage that they cannot
be routed on the Internet, so must be used only within a limited
scope.  In many situations, this is not a problem, but in certain
situations this can be problematic.  Though there is currently no
document that describes these situations, many of them are similar to
those described in [RFC6752].  See also
[I-D.ietf-v6ops-ula-usage-recommendations].

Not discussed in this document is the possibility of using the
technology described in [RFC6296] to work around some of the
limitations of PA GUAs and ULAs.

2.1.1.  Where to Use Addresses

As mentioned above, all interfaces in IPv6 always have a link-local
address.  This section addresses the question of when and where to
assign multi-hop addresses in addition to the LLA.  We consider four
options:

a.  Use only link-local addresses on all router interfaces.

b.  Assign multi-hop addresses to all link interfaces on each router,
    and use only a link-local address on the loopback interfaces.

c.  Assign multi-hop addresses to the loopback interface on each
    router, and use only a link-local address on all link interfaces.

   d.  Assign multi-hop addresses to both link and loopback interfaces
       on each router.

Option (a) means that the router cannot be reached (ping, management,
etc.) from farther than one-hop away.  The authors are not aware of
anyone using this option.

Option (b) means that the loopback interfaces are effectively
useless, since link-local addresses cannot be used for the purposes
that loopback interfaces are usually used for.  So option (b)
degenerates into option (d).

Thus the real choice comes down to option (c) vs. option (d).

Option (c) has two advantages over option (d).  The first advantage
is ease of configuration.  In a network with a large number of links,
the operator can just assign one multi-hop address to each router and
then enable the IGP, without going through the tedious process of
assigning and tracking the addresses on each link.  The second
advantage is security.  Since packets with link-local addresses
cannot be should not be routed, it is very difficult to attack the
associated nodes from an off-link device.  This implies less effort
around maintaining security ACLs.

Countering these advantages are various disadvantages to option (c)
compared with option (d):

o  It is not possible to ping a link-local-only interface from a
   device that is not directly attached to the link.  Thus, to
   troubleshoot, one must typically log into a device that is
   directly attached to the device in question, and execute the ping
   from there.

o  A traceroute passing over the link-local-only interface will
   return the loopback address of the router, rather than the address
   of the interface itself.

o  In cases of parallel point to point links it is difficult to
   determine which of the parallel links was taken when attempting to
   troubleshoot unless one sends packets directly between the two
   attached link-locals on the specific interfaces.  Since many
   network problems behave differently for traffic to/from a router
   than for traffic through the router(s) in question, this can pose
   a significant hurdle to some troubleshooting scenarios.

o  On some routers, by default the link-layer address of the
   interface is derived from the MAC address assigned to interface.
   When this is done, swapping out the interface hardware (e.g.

interface card) will cause the link-layer address to change.  In
some cases (peering config, ACLs, etc) this may require additional
changes.  However, many devices allow the link-layer address of an
interface to be explicitly configured, which avoids this issue.
This problem should fade away over time as more and more routers
select interface identifiers according to the rules in [RFC7217].

o  The practice of naming router interfaces using DNS names is
   difficult and not recommended when using link-locals only.  More
   generally, it is not recommended to put link-local addresses into
   DNS; see [RFC4472].

o  It is often not possible to identify the interface or link (in a
   database, email, etc) by giving just its address without also
   specifying the link in some manner.

It should be noted that it is quite possible for the same link-local
address to be assigned to multiple interfaces.  This can happen
because the MAC address is duplicated (due to manufacturing process
defaults or the use of virtualization), because a device deliberately
re-uses automatically-assigned link-local addresses on different
links, or because an operator manually assigns the same easy-to-type
link-local address to multiple interfaces.  All these are allowed in
IPv6 as long as the addresses are used on different links.

For more discussion on the pros and cons, see [RFC7404].  See also
[RFC5375] for IPv6 unicast address assignment considerations.

Today, most operators use option (d).

2.1.2.  Which Addresses to Use

Having considered above whether or not to use a "multi-hop address",
we now consider which of the addresses to use.

When selecting between these three "multi-hop address" types, one
needs to consider exactly how they will be used.  An important
consideration is how Internet traffic is carried across the core of
the network.  There are two main options: (1) the classic approach
where Internet traffic is carried as unlabeled traffic hop-by-hop
across the network, and (2) the more recent approach where Internet
traffic is carried inside an MPLS LSP (typically as part of a L3
VPN).

Under the classic approach:

o  PI GUAs are a very reasonable choice, if they are available.

   o  PA GUAs suffer from the "must renumber" and "difficult to multi-
      home" problems mentioned above.

   o  ULAs suffer from the "may be problematic" issues described above.

   Under the MPLS approach:

   o  PA GUAs are a reasonable choice, if they are available.

   o  PA GUAs suffer from the "must renumber" problem, but the
      "difficult to multi-home" problem does not apply.

   o  ULAs are a reasonable choice, since (unlike in the classic
      approach) these addresses are not visible to the Internet, so the
      problematic cases do not occur.

2.2.  Interfaces

2.2.1.  Mix IPv4 and IPv6 on the Same Layer-3 Interface?

   If a network is going to carry both IPv4 and IPv6 traffic, as many
   networks do today, then a question arises: Should an operator mix
   IPv4 and IPv6 traffic or keep them separated?  More specifically,
   should the design:

   a.  Mix IPv4 and IPv6 traffic on the same layer-3 interface, OR

   b.  Separate IPv4 and IPv6 by using separate interfaces (e.g., two
       physical links or two VLANs on the same link)?

   Option (a) implies a single layer-3 interface at each end of the
   connection with both IPv4 and IPv6 addresses; while option (b)
   implies two layer-3 interfaces at each end, one for IPv4 addresses
   and one with IPv6 addresses.

   The advantages of option (a) include:

   o  Requires only half as many layer 3 interfaces as option (b), thus
      providing better scaling;

   o  May require fewer physical ports, thus saving money and
      simplifying operations;

   o  Can make the QoS implementation much easier (for example, rate-
      limiting the combined IPv4 and IPv6 traffic to or from a
      customer);

   o  Works well in practice, as any increase in IPv6 traffic is usually
      counter-balanced by a corresponding decrease in IPv4 traffic to or
      from the same host (ignoring the common pattern of an overall
      increase in Internet usage);

   o  And is generally conceptually simpler.

   For these reasons, there is a relatively strong consensus in the
   operator community that option (a) is the preferred way to go.  Most
   networks today use option (a) wherever possible.

   However, there can be times when option (b) is the pragmatic choice.
   Most commonly, option (b) is used to work around limitations in
   network equipment.  One big example is the generally poor level of
   support today for individual statistics on IPv4 traffic vs IPv6
   traffic when option (a) is used.  Other, device-specific, limitations
   exist as well.  It is expected that these limitations will go away as
   support for IPv6 matures, making option (b) less and less attractive
   until the day that IPv4 is finally turned off.

## 2.3.  Static Routes

### 2.3.1.  Link-Local Next-Hop in a Static Route?

   For the most part, the use of static routes in IPv6 parallels their
   use in IPv4.  There is, however, one exception, which revolves around
   the choice of next-hop address in the static route.  Specifically,
   should an operator:

   a.  Use the far-end's link-local address as the next-hop address, OR

   b.  Use the far-end's GUA/ULA address as the next-hop address?

   Recall that the IPv6 specs for OSPF [RFC5340] and ISIS [RFC5308]
   dictate that they always use link-locals for next-hop addresses.  For
   static routes, [RFC4861] section 8 says:

      A router MUST be able to determine the link-local address for each
      of its neighboring routers in order to ensure that the target
      address in a Redirect message identifies the neighbor router by
      its link-local address.  For static routing, this requirement
      implies that the next-hop router's address should be specified
      using the link-local address of the router.

   This implies that using a GUA or ULA as the next hop will prevent a
   router from sending Redirect messages for packets that "hit" this
   static route.  All this argues for using a link-local as the next-hop
   address in a static route.

However, there are two cases where using a link-local address as the
next-hop clearly does not work.  One is when the static route is an
indirect (or multi-hop) static route.  The second is when the static
route is redistributed into another routing protocol.  In these
cases, the above text from RFC 4861 notwithstanding, either a GUA or
ULA must be used.

Furthermore, many network operators are concerned about the
dependency of the default link-local address on an underlying MAC
address, as described in the previous section.

Today most operators use GUAs as next-hop addresses.

## 2.4.  IGPs

### 2.4.1.  IGP Choice

One of the main decisions for a network operator looking to deploy
IPv6 is the choice of IGP (Interior Gateway Protocol) within the
network.  The main options are OSPF, IS-IS and EIGRP.  RIPng is
another option, but very few networks run RIP in the core these days,
so it is covered in a separate section below.

OSPF [RFC2328] [RFC5340] and IS-IS [RFC5120][RFC5120] are both
standardized link-state protocols.  Both protocols are widely
supported by vendors, and both are widely deployed.  By contrast,
EIGRP [RFC7868] is a Cisco proprietary distance-vector protocol.
EIGRP is rarely deployed in service-provider networks, but is quite
common in enterprise networks, which is why it is discussed here.

It is out of scope for this document to describe all the differences
between the three protocols; the interested reader can find books and
websites that go into the differences in quite a bit of detail.
Rather, this document simply highlights a few differences that can be
important to consider when designing IPv6 or dual-stack networks.

Versions: There are two versions of OSPF: OSPFv2 and OSPFv3.  The two
versions share many concepts, are configured in a similar manner and
seem very similar to most casual users, but have very different
packet formats and other "under the hood" differences.  The most
important difference is that OSPFv2 will only route IPv4, while
OSPFv3 will route both IPv4 and IPv6 (see [RFC5838]).  OSPFv2 was by
far the most widely deployed version of OSPF when this document was
published.  By contrast, both IS-IS and EIGRP have just a single
version, which can route both IPv4 and IPv6.

Transport.  IS-IS runs over layer 2 (e.g.  Ethernet).  This means
that the functioning of IS-IS has no dependencies on the IP layer: if

there is a problem at the IP layer (e.g. bad addresses), two routers
can still exchange IS-IS packets.  By contrast, OSPF and EIGRP both
run over the IP layer.  This means that the IP layer must be
configured and working OSPF or EIGRP packets to be exchanged between
routers.  For EIGRP, the dependency on the IP layer is simple: EIGRP
for IPv4 runs over IPv4, while EIGRP for IPv6 runs over IPv6.  For
OSPF, the story is more complex: OSPFv2 runs over IPv4, but OSPFv3
can run over either IPv4 or IPv6.  Thus it is possible to route both
IPv4 and IPv6 with OSPFv3 running over IPv6 or with OSPFv3 running
over IPv4.  This means that there are number of choices for how to
run OSPF in a dual-stack network:

o  Use OSPFv2 for routing IPv4 , and OSPFv3 running over IPv6 for
   routing IPv6, OR

o  Use OSPFv3 running over IPv6 for routing both IPv4 and IPv6, OR

o  Use OSPFv3 running over IPv4 for routing both IPv4 and IPv6.


Summarization and MPLS: For most casual users, the three protocols
are fairly similar in what they can do, with two glaring exceptions:
summarization and MPLS.  For summarization, both OSPF and IS-IS have
the concept of summarization between areas, but the two area concepts
are quite different, and an area design that works for one protocol
will usually not work for the other.  EIGRP has no area concept, but
has the ability to summarize at any router.  Thus a large network
will typically have a very different OSPF, IS-IS and EIGRP designs,
which is important to keep in mind if you are planning on using one
protocol to route IPv4 and a different protocol for IPv6.  The other
difference is that OSPF and IS-IS both support RSVP-TE, a widely-used
MPLS signaling protocol, while EIGRP does not: this is due to OSPF
and IS-IS both being link-state protocols while EIGRP is a distance-
vector protocol.

The table below sets out possible combinations of protocols to route
both IPv4 and IPv6, and makes some observations on each combination.
Here "EIGRP-v4" means "EIGRP for IPv4" and similarly for "EIGRP-v6".
For OSPFv3, it is possible to run it over either IPv4 or IPv6; this
is not indicated in the table.

| IGP for IPv4 | IGP for IPv6 | Protocol separation | Similar configuration possible | Multiple Known Deployments |
|------|------|------|------|------|
| | | | | |
| OSPFv2 | OSPFv3 | YES | YES | YES (8) |
| OSPFv2 | IS-IS | YES | - | YES (3) |
| OSPFv2 | EIGRP-v6 | YES | - | - |
| OSPFv3 | OSPFv3 | NO | YES | - |
| OSPFv3 | IS-IS | YES | - | - |
| OSPFv3 | EIGRP-v6 | YES | - | - |
| IS-IS | OSPFv3 | YES | - | YES (2) |
| IS-IS | IS-IS | - | YES | YES (12) |
| IS-IS | EIGRP-v6 | YES | - | - |
| EIGRP-v4 | OSPFv3 | YES | - | ? (1) |
| EIGRP-v4 | IS-IS | YES | - | - |
| EIGRP-v4 | EIGRP-v6 | - | YES | ? (2) |

In the column "Multiple Known Deployments", a YES indicates that a
significant number of production networks run this combination, with
the number of such networks indicated in parentheses following, while
a "?" indicates that the authors are only aware of one or two small
networks that run this combination.  Data for this column was
gathered from an informal poll of operators on a number of mailing
lists.  This poll was not intended to be a thorough scientific study
of IGP choices, but to provide a snapshot of known operator choices
at the time of writing (Mid-2015) for successful production dual
stack network deployments.  There were twenty six (26) network
implementations represented by 17 respondents.  Some respondents
provided information on more then one network or network deployment.
Due to privacy considerations, the networks' represented and
respondents are not listed in this document.

A number of combinations are marked as offering "Protocol separation".  These options use a different IGP protocol for IPv4 vs IPv6.  With these options, a problem with routing IPv6 is unlikely to affect IPv4 or visa-versa.  Some operator may consider this as a benefit when first introducing dual stack capabilities or for ongoing technical reasons.

Three combinations are marked "Similar configuration possible".  This means it is possible (but not required) to use very similar IGP configuration for IPv4 and IPv6: for example, the same area boundaries, area numbering, link costing, etc.  If you are happy with your IPv4 IGP design, then this will likely be a consideration.  By contrast, the options that use, for example, IS-IS for one IP version and OSPF for the other version will require considerably different configuration, and will also require the operations staff to become familiar with the difference between the two protocols.

It should be noted that a number of ISPs have run OSPF as their IPv4 IGP for quite a few years, but have selected IS-IS as their IPv6 IGP.  However, there are very few (none?) that have made the reverse choice.  This is, in part, because routers generally support more nodes in an IS-IS area than in the corresponding OSPF area, and because IS-IS is seen as more secure because it runs at layer 2.

2.4.2.  IS-IS Topology Mode

When IS-IS is used to route both IPv4 and IPv6, then there is an additional choice of whether to run IS-IS in single-topology or multi-topology mode.

With single-topology mode (also known as Native mode) [RFC5308]:

o  IS-IS keeps a single link-state database for both IPv4 and IPv6.

o  There is a single set of link costs which apply to both IPv4 and IPv6.

o  All links in the network must support both IPv4 and IPv6, as the calculation of routes does not take this into account.  If some links do not support IPv6 (or IPv4), then packets may get routed across links where support is lacking and get dropped.  This can cause problems if some network devices do not support IPv6 (or IPv4).

o  It is also important to keep the previous point in mind when adding or removing support for either IPv4 or IPv6.

With multi-topology mode [RFC5120]:

   o  IS-IS keeps two link-state databases, one for IPv4 and one for
      IPv6.

   o  IPv4 and IPv6 can have separate link metrics.  Note that most
      implementations today require separate link metrics: a number of
      operators have rudely discovered that they have forgotten to
      configure the IPv6 metric until sometime after deploying IPv6 in
      multi-topology mode!

   o  Some links can be IPv4-only, some IPv6-only, and some dual-stack.
      Routes to IPv4 and IPv6 addresses are computed separately and may
      take different paths even if the addresses are located on the same
      remote device.

   o  The previous point may help when adding or removing support for
      either IPv4 or IPv6.

   In the informal poll of operators, out of 12 production networks that
   ran IS-IS for both IPv4 and IPv6, 6 used single topology mode, 4 used
   multi-topology mode, and 2 did not specify.  One motivation often
   cited by then operators for using Single Topology mode was because
   some device did not support multi-topology mode.

   When asked, many people feel multi-topology mode is superior to
   single-topology mode because it provides greater flexibility at
   minimal extra cost.  Never-the-less, as shown by the poll results, a
   number of operators have used single-topology mode successfully.

   Note that this issue does not come up with OSPF, since there is
   nothing that corresponds to IS-IS single-topology mode with OSPF.

2.4.3.  RIP / RIPng

   A protocol option not described in the table above is RIP for IPv4
   and RIPng for IPv6 [RFC2080].  These are distance vector protocols
   that are almost universally considered to be inferior to OSPF, IS-IS,
   or EIGRP for general use.

   However, there is one specialized use where RIP/RIPng is still
   considered to be appropriate: in star topology networks where a
   single core device has lots and lots of links to edge devices and
   each edge device has only a single path back to the core.  In such
   networks, the single path means that the limitations of RIP/RIPng are
   mostly not relevant and the very light-weight nature of RIP/RIPng
   gives it an advantage over the other protocols mentioned above.  One
   concrete example of this scenario is the use of RIP/RIPng between
   cable modems and the CMTS.

2.5.  BGP

2.5.1.  Which Transport for Which Routes?

   BGP these days is multi-protocol.  It can carry routes of many
   different types, or more precisely, many different AFI/SAFI
   combinations.  It can also carry routes when the BGP session, or more
   accurately the underlying TCP connection, runs over either IPv4 or
   IPv6 (here referred to as either "IPv4 transport" or "IPv6
   transport").  Given this flexibility, one of the biggest questions
   when deploying BGP in a dual-stack network is the question of which
   route types should be carried over sessions using IPv4 transport and
   which should be carried over sessions using IPv6 transport.

   This section discusses this question for the three most-commonly-used
   SAFI values: unlabeled (SAFI 1), labeled (SAFI 4) and VPN (SAFI 128).
   Though we do not explicitly discuss other SAFI values, many of the
   comments here can be applied to the other values.

   Consider the following table:

```
+---------------+----------+----------------------------+
| Route Family  | Transport | Comments                   |
+---------------+----------+----------------------------+
|               |          |                            |
+---------------+----------+----------------------------+
| Unlabeled IPv4 |   IPv4  | Works well                 |
+---------------+----------+----------------------------+
| Unlabeled IPv4 |   IPv6  | Next-hop                   |
+---------------+----------+----------------------------+
| Unlabeled IPv6 |   IPv4  | Next-hop                   |
+---------------+----------+----------------------------+
| Unlabeled IPv6 |   IPv6  | Works well                 |
+---------------+----------+----------------------------+
|               |          |                            |
+---------------+----------+----------------------------+
|  Labeled IPv4 |   IPv4   | Works well                 |
+---------------+----------+----------------------------+
|  Labeled IPv4 |   IPv6   | Next-hop                   |
+---------------+----------+----------------------------+
|  Labeled IPv6 |   IPv4   | (6PE) Works well           |
+---------------+----------+----------------------------+
|  Labeled IPv6 |   IPv6   | Next-hop or MPLS over IPv6 |
+---------------+----------+----------------------------+
|               |          |                            |
+---------------+----------+----------------------------+
|    VPN IPv4   |   IPv4   | Works well                 |
+---------------+----------+----------------------------+
|    VPN IPv4   |   IPv6   | Next-hop                   |
+---------------+----------+----------------------------+
|    VPN IPv6   |   IPv4   | (6VPE) Works well          |
+---------------+----------+----------------------------+
|    VPN IPv6   |   IPv6   | Next-hop or MPLS over IPv6 |
+---------------+----------+----------------------------+
```

   The first column in this table lists various route families, where
   "unlabeled" means SAFI 1, "labeled" means the routes carry an MPLS
   label (SAFI 4, see [RFC3107]), and "VPN" means the routes are
   normally associated with a layer-3 VPN (SAFI 128, see [RFC4364]).
   The second column lists the protocol used to transport the BGP
   session, frequently specified by giving either an IPv4 or IPv6
   address in the "neighbor" statement.

   The third column comments on the combination in the first two
   columns:

   o  For combinations marked "Works well", these combinations are
      standardized, widely supported and widely deployed.

   o  For combinations marked "Next-hop", these combinations are not
      standardized and are less-widely supported.  These combinations
      all have the "next-hop mismatch" problem: the transported route
      needs a next-hop address from the other address family than the
      transport address (for example, an IPv4 route needs an IPv4 next-
      hop, even when transported over IPv6).  Some vendors have
      implemented ways to solve this problem for specific combinations,
      but for combinations marked "next-hop", these solutions have not
      been standardized (cf. 6PE and 6VPE, where the solution has been
      standardized).

   o  For combinations marked as "Next-hop or MPLS over IPv6", these
      combinations either require a non-standard solution to the next-
      hop problem, or require MPLS over IPv6.  At the time of writing,
      MPLS over IPv6 is not widely supported or deployed.

   Also, it is important to note that changing the set of address
   families being carried over a BGP session requires the BGP session to
   be reset (unless something like [I-D.ietf-idr-dynamic-cap] or
   [I-D.ietf-idr-bgp-multisession] is in use).  This is generally more
   of an issue with eBGP sessions than iBGP sessions: for iBGP sessions
   it is common practice for a router to have two iBGP sessions, one to
   each member of a route reflector pair, so one can change the set of
   address families on first one of the sessions and then the other.

   The following subsections discuss specific combinations in more
   detail.

2.5.1.1.  BGP Sessions for Unlabeled Routes

   Unlabeled routes are commonly carried on eBGP sessions, as well as on
   iBGP sessions in networks where Internet traffic is carried unlabeled
   across the network.

   In these scenarios, there are three reasonable choices:

   a.  Carry unlabeled IPv4 and IPv6 routes over IPv4, OR

   b.  Carry unlabeled IPv4 and IPv6 routes over IPv6, OR

   c.  Carry unlabeled IPv4 routes over IPv4, and unlabeled IPv6 routes
       over IPv6

   Options (a) and (b) have the advantage that one one BGP session is
   required between pairs of routers.  However, option (c) is widely
   considered to be the best choice.  There are several reasons for this
   :

   o  It gives a clean separation between IPv4 and IPv6.  This can be
      especially useful when first deploying IPv6 and troubleshooting
      resulting problems.

   o  This avoids the next-hop problem described above.

   o  The status of the routes follows the status of the underlying
      transport.  If, for example, the IPv6 data path between the two
      BGP speakers fails, then the IPv6 session between the two speakers
      will fail and the IPv6 routes will be withdrawn, which will allow
      the traffic to be re-routed elsewhere.  By contrast, if the IPv6
      routes were transported over IPv4, then the failure of the IPv6
      data path might leave a working IPv4 data path, so the BGP session
      would remain up and the IPv6 routes would not be withdrawn, and
      thus the IPv6 traffic would be sent into a black hole.

   o  It avoids resetting the BGP session when adding IPv6 to an
      existing session, or when removing IPv4 from an existing session.

   Rarely, there are situations where option (c) is not practical.  In
   those cases today, most operators use option (a), carrying both route
   types over a single BGP session.

2.5.1.2.  BGP sessions for Labeled or VPN Routes

   When carrying labeled or VPN routes, the only widely-supported
   solution at time of writing is to carry both route types over IPv4.
   This may change in as MPLS over IPv6 becomes more widely implemented.

   There are two options when carrying both over IPv4:

   a.  Carry all routes over a single BGP session, OR

   b.  Carry the routes over multiple BGP sessions (e.g. one for VPN
       IPv4 routes and one for VPN IPv6 routes)

   Using a single session is usually simplest for an iBGP session going
   to a route reflector handling both route families.  Using a single
   session here usually means that the BGP session will reset when
   changing the set of address families, but as noted above, this is
   usually not a problem when redundant route reflectors are involved.

   In eBGP situations, two sessions are usually more appropriate.
   [JUSTIFICATION?]

2.5.2.  eBGP Endpoints: Global or Link-Local Addresses?

   When running eBGP over IPv6, there are two options for the addresses
   to use at each end of the eBGP session (or more properly, the
   underlying TCP session):

   a.  Use link-local addresses for the eBGP session, OR

   b.  Use global addresses for the eBGP session.

   Note that the choice here is the addresses to use for the eBGP
   sessions, and not whether the link itself has global (or unique-
   local) addresses.  In particular, it is quite possible for the eBGP
   session to use link-local addresses even when the link has global
   addresses.

   The big attraction for option (a) is security: an eBGP session using
   link-local addresses is extremely difficult to attack from a device
   that is off-link.  This provides very strong protection against TCP
   RST and similar attacks.  Though there are other ways to get an
   equivalent level of security (e.g.  GTSM [RFC5082], MD5 [RFC5925], or
   ACLs), these other ways require additional configuration which can be
   forgotten or potentially mis-configured.

   However, there are a number of small disadvantages to using link-
   local addresses:

   o  Using link-local addresses only works for single-hop eBGP
      sessions; it does not work for multi-hop sessions.

   o  One must use "next-hop self" at both endpoints, otherwise re-
      advertising routes learned via eBGP into iBGP will not work.
      (Some products enable "next-hop self" in this situation
      automatically).

   o  Operators and their tools are used to referring to eBGP sessions
      by address only, something that is not possible with link-local
      addresses.

   o  If one is configuring parallel eBGP sessions for IPv4 and IPv6
      routes, then using link-local addresses for the IPv6 session
      introduces extra operational differences between the two sessions
      which could otherwise be avoided.

   o  On some products, an eBGP session using a link-local address is
      more complex to configure than a session that uses a global
      address.

o  If hardware or other issues cause one to move the cable to a
   different local interface, then reconfiguration is required at
   both ends: at the local end because the interface has changed (and
   with link-local addresses, the interface must always be specified
   along with the address), and at the remote end because the link-
   local address has likely changed.  (Contrast this with using
   global addresses, where less re-configuration is required at the
   local end, and no reconfiguration is required at the remote end).

o  Finally, a strict application of [RFC2545] forbids running eBGP
   between link-local addresses, as [RFC2545] requires the BGP next-
   hop field to contain at least a global address.

For these reasons, most operators today choose to have their eBGP
sessions use global addresses.

3.  General Observations

There are two themes that run though many of the design choices in
this document.  This section presents some general discussion on
these two themes.

3.1.  Use of Link-Local Addresses

The proper use of link-local addresses is a common theme in the IPv6
network design choices.  Link-layer addresses are, of course, always
present in an IPv6 network, but current network design practice
mostly ignores them, despite efforts such as [RFC7404].

There are three main reasons for this current practice:

o  Network operators are concerned about the volatility of link-local
   addresses based on MAC addresses, despite the fact that this
   concern can be overcome by manually-configuring link-local
   addresses;

o  It is very difficult to impossible to ping a link-local address
   from a device that is not on the same subnet.  This is a
   troubleshooting disadvantage, though it can also be viewed as a
   security advantage.

o  Most operators are currently running networks that carry both IPv4
   and IPv6 traffic, and wish to harmonize their IPv4 and IPv6 design
   and operational practices where possible.

3.2.  Separation of IPv4 and IPv6

   Currently, most operators are running or planning to run networks
   that carry both IPv4 and IPv6 traffic.  Hence the question: To what
   degree should IPv4 and IPv6 be kept separate?  As can be seen above,
   this breaks into two sub-questions: To what degree should IPv4 and
   IPv6 traffic be kept separate, and to what degree should IPv4 and
   IPv6 routing information be kept separate?

   The general consensus around the first question is that IPv4 and IPv6
   traffic should generally be mixed together.  This recommendation is
   driven by the operational simplicity of mixing the traffic, plus the
   general observation that the service being offered to the end user is
   Internet connectivity and most users do not know or care about the
   differences between IPv4 and IPv6.  Thus it is very desirable to mix
   IPv4 and IPv6 on the same link to the end user.  On other links,
   separation is possible but more operationally complex, though it does
   occasionally allow the operator to work around limitations on network
   devices.  The situation here is roughly comparable to IP and MPLS
   traffic: many networks mix the two traffic types on the same links
   without issues.

   By contrast, there is more of an argument for carrying IPv6 routing
   information over IPv6 transport, while leaving IPv4 routing
   information on IPv4 transport.  By doing this, one gets fate-sharing
   between the control and data plane for each IP protocol version: if
   the data plane fails for some reason, then often the control plane
   will too.

4.  IANA Considerations

   This document makes no requests of IANA.

5.  Security Considerations

   This document introduces no new security considerations that are not
   already documented elsewhere.

   The following is a brief list of pointers to documents related to the
   topics covered above that the reader may wish to review for security
   considerations.

   For general IPv6 security, [RFC4942] provides guidance on security
   considerations around IPv6 transition and coexistence.

   For OSPFv3, the base protocol specification [RFC5340] has a short
   security considerations section which notes that the fundamental

mechanism for protecting OSPFv3 from attacks is the mechanism
described in [RFC4552].

For IS-IS, [RFC5308] notes that ISIS for IPv6 raises no new security
considerations over ISIS for IPv4 over those documented in [ISO10589]
and [RFC5304].

For BGP, [RFC2545] notes that BGP for IPv6 raises no new security
considerations over those present in BGP for IPv4.  However, there
has been much discussion of BGP security recently, and the interested
reader is referred to the documents of the IETF's SIDR working group.

6.  Acknowledgements

Many, many people in the V6OPS working group provided comments and
suggestions that made their way into this document.  A partial list
includes: Rajiv Asati, Fred Baker, Michael Behringer, Marc Blanchet,
Ron Bonica, Randy Bush, Cameron Byrne, Brian Carpenter, KK
Chittimaneni, Tim Chown, Lorenzo Colitti, Gert Doering, Francis
Dupont, Bill Fenner, Kedar K Gaonkar, Chris Grundemann, Steinar Haug,
Ray Hunter, Joel Jaeggli, Victor Kuarsingh, Jen Linkova, Ivan
Pepelnjak, Alexandru Petrescu, Rob Shakir, Mark Smith, Jean-Francois
Tremblay, Dave Thaler, Tina Tsou, Eric Vyncke, Dan York, and
Xuxiaohu.

The authors would also like to thank Pradeep Jain and Alastair
Johnson for helpful comments on a very preliminary version of this
document.

7.  Informative References

   [I-D.ietf-idr-bgp-multisession]
           Scudder, J., Appanna, C., and I. Varlashkin, "Multisession
           BGP", draft-ietf-idr-bgp-multisession-07 (work in
           progress), September 2012.

   [I-D.ietf-idr-dynamic-cap]
           Ramachandra, S. and E. Chen, "Dynamic Capability for BGP-
           4", draft-ietf-idr-dynamic-cap-14 (work in progress),
           December 2011.

   [I-D.ietf-v6ops-ula-usage-recommendations]
           Liu, B. and S. Jiang, "Considerations For Using Unique
           Local Addresses", draft-ietf-v6ops-ula-usage-
           recommendations-05 (work in progress), May 2015.

[ISO10589]
          International Standards Organization, "Intermediate system
          to Intermediate system intra-domain routeing information
          exchange protocol for use in conjunction with the protocol
          for providing the connectionless-mode Network Service (ISO
          8473)", International Standard 10589:2002, Nov 2002.

[RFC2080]  Malkin, G. and R. Minnear, "RIPng for IPv6", RFC 2080,
          DOI 10.17487/RFC2080, January 1997,
          <http://www.rfc-editor.org/info/rfc2080>.

[RFC2328]  Moy, J., "OSPF Version 2", STD 54, RFC 2328,
          DOI 10.17487/RFC2328, April 1998,
          <http://www.rfc-editor.org/info/rfc2328>.

[RFC2545]  Marques, P. and F. Dupont, "Use of BGP-4 Multiprotocol
          Extensions for IPv6 Inter-Domain Routing", RFC 2545,
          DOI 10.17487/RFC2545, March 1999,
          <http://www.rfc-editor.org/info/rfc2545>.

[RFC3107]  Rekhter, Y. and E. Rosen, "Carrying Label Information in
          BGP-4", RFC 3107, DOI 10.17487/RFC3107, May 2001,
          <http://www.rfc-editor.org/info/rfc3107>.

[RFC4193]  Hinden, R. and B. Haberman, "Unique Local IPv6 Unicast
          Addresses", RFC 4193, DOI 10.17487/RFC4193, October 2005,
          <http://www.rfc-editor.org/info/rfc4193>.

[RFC4291]  Hinden, R. and S. Deering, "IP Version 6 Addressing
          Architecture", RFC 4291, DOI 10.17487/RFC4291, February
          2006, <http://www.rfc-editor.org/info/rfc4291>.

[RFC4364]  Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private
          Networks (VPNs)", RFC 4364, DOI 10.17487/RFC4364, February
          2006, <http://www.rfc-editor.org/info/rfc4364>.

[RFC4472]  Durand, A., Ihren, J., and P. Savola, "Operational
          Considerations and Issues with IPv6 DNS", RFC 4472,
          DOI 10.17487/RFC4472, April 2006,
          <http://www.rfc-editor.org/info/rfc4472>.

[RFC4552]  Gupta, M. and N. Melam, "Authentication/Confidentiality
          for OSPFv3", RFC 4552, DOI 10.17487/RFC4552, June 2006,
          <http://www.rfc-editor.org/info/rfc4552>.

   [RFC4852]  Bound, J., Pouffary, Y., Klynsma, S., Chown, T., and D.
              Green, "IPv6 Enterprise Network Analysis - IP Layer 3
              Focus", RFC 4852, DOI 10.17487/RFC4852, April 2007,
              <http://www.rfc-editor.org/info/rfc4852>.

   [RFC4861]  Narten, T., Nordmark, E., Simpson, W., and H. Soliman,
              "Neighbor Discovery for IP version 6 (IPv6)", RFC 4861,
              DOI 10.17487/RFC4861, September 2007,
              <http://www.rfc-editor.org/info/rfc4861>.

   [RFC4942]  Davies, E., Krishnan, S., and P. Savola, "IPv6 Transition/
              Co-existence Security Considerations", RFC 4942,
              DOI 10.17487/RFC4942, September 2007,
              <http://www.rfc-editor.org/info/rfc4942>.

   [RFC5082]  Gill, V., Heasley, J., Meyer, D., Savola, P., Ed., and C.
              Pignataro, "The Generalized TTL Security Mechanism
              (GTSM)", RFC 5082, DOI 10.17487/RFC5082, October 2007,
              <http://www.rfc-editor.org/info/rfc5082>.

   [RFC5120]  Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi
              Topology (MT) Routing in Intermediate System to
              Intermediate Systems (IS-ISs)", RFC 5120,
              DOI 10.17487/RFC5120, February 2008,
              <http://www.rfc-editor.org/info/rfc5120>.

   [RFC5220]  Matsumoto, A., Fujisaki, T., Hiromi, R., and K. Kanayama,
              "Problem Statement for Default Address Selection in Multi-
              Prefix Environments: Operational Issues of RFC 3484
              Default Rules", RFC 5220, DOI 10.17487/RFC5220, July 2008,
              <http://www.rfc-editor.org/info/rfc5220>.

   [RFC5304]  Li, T. and R. Atkinson, "IS-IS Cryptographic
              Authentication", RFC 5304, DOI 10.17487/RFC5304, October
              2008, <http://www.rfc-editor.org/info/rfc5304>.

   [RFC5308]  Hopps, C., "Routing IPv6 with IS-IS", RFC 5308,
              DOI 10.17487/RFC5308, October 2008,
              <http://www.rfc-editor.org/info/rfc5308>.

   [RFC5340]  Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF
              for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008,
              <http://www.rfc-editor.org/info/rfc5340>.

   [RFC5375]  Van de Velde, G., Popoviciu, C., Chown, T., Bonness, O.,
              and C. Hahn, "IPv6 Unicast Address Assignment
              Considerations", RFC 5375, DOI 10.17487/RFC5375, December
              2008, <http://www.rfc-editor.org/info/rfc5375>.

   [RFC5838]  Lindem, A., Ed., Mirtorabi, S., Roy, A., Barnes, M., and
              R. Aggarwal, "Support of Address Families in OSPFv3",
              RFC 5838, DOI 10.17487/RFC5838, April 2010,
              <http://www.rfc-editor.org/info/rfc5838>.

   [RFC5887]  Carpenter, B., Atkinson, R., and H. Flinck, "Renumbering
              Still Needs Work", RFC 5887, DOI 10.17487/RFC5887, May
              2010, <http://www.rfc-editor.org/info/rfc5887>.

   [RFC5925]  Touch, J., Mankin, A., and R. Bonica, "The TCP
              Authentication Option", RFC 5925, DOI 10.17487/RFC5925,
              June 2010, <http://www.rfc-editor.org/info/rfc5925>.

   [RFC5963]  Gagliano, R., "IPv6 Deployment in Internet Exchange Points
              (IXPs)", RFC 5963, DOI 10.17487/RFC5963, August 2010,
              <http://www.rfc-editor.org/info/rfc5963>.

   [RFC6180]  Arkko, J. and F. Baker, "Guidelines for Using IPv6
              Transition Mechanisms during IPv6 Deployment", RFC 6180,
              DOI 10.17487/RFC6180, May 2011,
              <http://www.rfc-editor.org/info/rfc6180>.

   [RFC6296]  Wasserman, M. and F. Baker, "IPv6-to-IPv6 Network Prefix
              Translation", RFC 6296, DOI 10.17487/RFC6296, June 2011,
              <http://www.rfc-editor.org/info/rfc6296>.

   [RFC6342]  Koodli, R., "Mobile Networks Considerations for IPv6
              Deployment", RFC 6342, DOI 10.17487/RFC6342, August 2011,
              <http://www.rfc-editor.org/info/rfc6342>.

   [RFC6752]  Kirkham, A., "Issues with Private IP Addressing in the
              Internet", RFC 6752, DOI 10.17487/RFC6752, September 2012,
              <http://www.rfc-editor.org/info/rfc6752>.

   [RFC6782]  Kuarsingh, V., Ed. and L. Howard, "Wireline Incremental
              IPv6", RFC 6782, DOI 10.17487/RFC6782, November 2012,
              <http://www.rfc-editor.org/info/rfc6782>.

   [RFC6879]  Jiang, S., Liu, B., and B. Carpenter, "IPv6 Enterprise
              Network Renumbering Scenarios, Considerations, and
              Methods", RFC 6879, DOI 10.17487/RFC6879, February 2013,
              <http://www.rfc-editor.org/info/rfc6879>.

   [RFC6883]  Carpenter, B. and S. Jiang, "IPv6 Guidance for Internet
              Content Providers and Application Service Providers",
              RFC 6883, DOI 10.17487/RFC6883, March 2013,
              <http://www.rfc-editor.org/info/rfc6883>.

   [RFC7010]  Liu, B., Jiang, S., Carpenter, B., Venaas, S., and W.
              George, "IPv6 Site Renumbering Gap Analysis", RFC 7010,
              DOI 10.17487/RFC7010, September 2013,
              <http://www.rfc-editor.org/info/rfc7010>.

   [RFC7217]  Gont, F., "A Method for Generating Semantically Opaque
              Interface Identifiers with IPv6 Stateless Address
              Autoconfiguration (SLAAC)", RFC 7217,
              DOI 10.17487/RFC7217, April 2014,
              <http://www.rfc-editor.org/info/rfc7217>.

   [RFC7381]  Chittimaneni, K., Chown, T., Howard, L., Kuarsingh, V.,
              Pouffary, Y., and E. Vyncke, "Enterprise IPv6 Deployment
              Guidelines", RFC 7381, DOI 10.17487/RFC7381, October 2014,
              <http://www.rfc-editor.org/info/rfc7381>.

   [RFC7404]  Behringer, M. and E. Vyncke, "Using Only Link-Local
              Addressing inside an IPv6 Network", RFC 7404,
              DOI 10.17487/RFC7404, November 2014,
              <http://www.rfc-editor.org/info/rfc7404>.

   [RFC7868]  Savage, D., Ng, J., Moore, S., Slice, D., Paluch, P., and
              R. White, "Cisco's Enhanced Interior Gateway Routing
              Protocol (EIGRP)", RFC 7868, DOI 10.17487/RFC7868, May
              2016, <http://www.rfc-editor.org/info/rfc7868>.

   [v6-addressing-plan]
              SurfNet, "Preparing an IPv6 Address Plan", 2013,
              <http://www.ripe.net/lir-services/training/material/
              IPv6-for-LIRs-Training-Course/
              Preparing-an-IPv6-Addressing-Plan.pdf>.

Authors' Addresses

   Philip Matthews
   Nokia
   600 March Road
   Ottawa, Ontario  K2K 2E6
   Canada

   Phone: +1 613-784-3139
   Email: philip_matthews@magma.ca

Victor Kuarsingh
Cisco
88 Queens Quay
Toronto, ON  M5J0B8
Canada

Email: victor@jvknet.com