Network Working Group                                          K. Patel
Internet-Draft                                             Arrcus, Inc.
Intended status: Standards Track                             A. Lindem
Expires: September 4, 2018                               Cisco Systems
                                                             S. Zandi
                                                             Linkedin
                                                        W. Henderickx
                                                                Nokia
                                                        March 3, 2018

             Shortest Path Routing Extensions for BGP Protocol
                     draft-keyupate-lsvr-bgp-spf-00.txt

Abstract

   Many Massively Scaled Data Centers (MSDCs) have converged on
   simplified layer 3 routing.  Furthermore, requirements for
   operational simplicity have lead many of these MSDCs to converge on
   BGP as their single routing protocol for both their fabric routing
   and their Data Center Interconnect (DCI) routing.  This document
   describes a solution which leverages BGP Link-State distribution and
   the Shortest Path First algorithm similar to Internal Gateway
   Protocols (IGPs) such as OSPF.

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 September 4, 2018.

Table of Contents

## 1.  Introduction

   Many Massively Scaled Data Centers (MSDCs) have converged on
   simplified layer 3 routing.  Furthermore, requirements for
   operational simplicity have lead many of these MSDCs to converge on
   BGP [RFC4271] as their single routing protocol for both their fabric
   routing and their Data Center Interconnect (DCI) routing.
   Requirements and procedures for using BGP are described in [RFC7938].
   This document describes an alternative solution which leverages BGP-
   LS [RFC7752] and the Shortest Path First algorithm similar to
   Internal Gateway Protocols (IGPs) such as OSPF [RFC2328].

   [RFC4271] defines the Decision Process that is used to select routes
   for subsequent advertisement by applying the policies in the local
   Policy Information Base (PIB) to the routes stored in its Adj-RIBs-
   In.  The output of the Decision Process is the set of routes that are
   announced by a BGP speaker to its peers.  These selected routes are
   stored by a BGP speaker in the speaker's Adj-RIBs-Out according to
   policy.

   [RFC7752] describes a mechanism by which link-state and TE
   information can be collected from networks and shared with external
   components using BGP.  This is achieved by defining NLRI carried
   within BGP-LS AFI and BGP-LS SAFIs.  The BGP-LS extensions defined in
   [RFC7752] makes use of the Decision Process defined in [RFC4271].

   This document augments [RFC7752] by replacing its use of the existing
   Decision Process.  The BGP-LS-SPF and BGP-LS-SPF-VPN AFI/SAFI are
   introduced to insure backward compatibility.  The Phase 1 and 2
   decision functions of the Decision Process are replaced with the
   Shortest Path Algorithm (SPF) also known as the Dijkstra Algorithm.
   The Phase 3 decision function is also simplified since it is no
   longer dependent on the previous phases.  This solution avails the
   benefits of both BGP and SPF-based IGPs.  These include TCP based
   flow-control, no periodic link-state refresh, and completely
   incremental NLRI advertisement.  These advantages can reduce the
   overhead in MSDCs where there is a high degree of Equal Cost Multi-
   Path (ECMPs) and the topology is very stable.  Additionally, using a
   SPF-based computation can support fast convergence and the
   computation of Loop-Free Alternatives (LFAs) [RFC5286] in the event
   of link failures.  Furthermore, a BGP based solution lends itself to
   multiple peering models including those incorporating route-
   reflectors [RFC4456] or controllers.

Support for Multiple Topology Routing (MTR) as described in [RFC4915]
is an area for further study dependent on deployment requirements.

## 1.1.  BGP Shortest Path First (SPF) Motivation

Given that [RFC7938] already describes how BGP could be used as the
sole routing protocol in an MSDC, one might question the motivation
for defining an alternate BGP deployment model when a mature solution
exists.  For both alternatives, BGP offers the operational benefits
of a single routing protocol.  However, BGP SPF offers some unique
advantages above and beyond standard BGP distance-vector routing.

A primary advantage is that all BGP speakers in the BGP SPF routing
domain will have a complete view of the topology.  This will allow
support of ECMP, IP fast-reroute (e.g., Loop-Free Alternatives),
Shared Risk Link Groups (SRLGs), and other routing enhancements
without advertisement of addition BGP paths or other extensions.  In
short, the advantages of an IGP such as OSPF [RFC2328] are availed in
BGP.

With the simplified BGP decision process as defined in Section 5.1,
NLRI changes can be disseminated throughout the BGP routing domain
much more rapidly (equivalent to IGPs with the proper
implementation).

Another primary advantage is a potential reduction in NLRI
advertisement.  With standard BGP distance-vector routing, a single
link failure may impact 100s or 1000s prefixes and result in the
withdrawal or re-advertisement of the attendant NLRI.  With BGP SPF,
only the BGP speakers corresponding to the link NLRI need withdraw
the corresponding BGP-LS Link NLRI.  This advantage will contribute
to both faster convergence and better scaling.

With controller and route-reflector peering models, BGP SPF
advertisement and distributed computation require a minimal number of
sessions and copies of the NLRI since only the latest verion of the
NLRI from the originator is required.  Given that verification of the
adjacencies is done outside of BGP (see Section 2), each BGP speaker
will only need as many sessions and copies of the NLRI as required
for redundancy (e.g., one for SPF computation and another for
backup).  Functions such as Optimized Route Reflection (ORR) are
supported without extension by virture of the primary advantages.
Additionally, a controller could inject topology that is learned
outside the BGP routing domain.

Given that controllers are already consuming BGP-LS NLRI [RFC7752],
reusing for the BGP-LS SPF leverages the existing controller
implementations.

Another potential advantage of BGP SPF is that both IPv6 and IPv4 can
be supported in the same address family using the same topology.
Although not described in this version of the document, multi-
topology extensions can be used to support separate IPv4, IPv6,
unicast, and multicast topologies while sharing the same NLRI.

Finally, the BGP SPF topology can be used as an underlay for other
BGP address families (using the existing model) and realize all the
above advantages.  A simplified peering model using IPv6 link-local
addresses as next-hops can be deployed similar to [RFC5549].

## 1.2.  Requirements Language

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 RFC 2119 [RFC2119].

## 2.  BGP Peering Models

Depending on the requirements, scaling, and capabilities of the BGP
speakers, various peering models are supported.  The only requirement
is that all BGP speakers in the BGP SPF routing domain receive link-
state NLRI on a timely basis, run an SPF calculation, and update
their data plane appropriately.  The content of the Link NLRI is
described in Section 4.2.

## 2.1.  BGP Single-Hop Peering on Network Node Connections

The simplest peering model is the one described in section 5.2.1 of
[RFC7938].  In this model, EBGP single-hop sessions are established
over direct point-to-point links interconnecting the network nodes.
For the purposes of BGP SPF, Link NLRI is only advertised if a
single-hop BGP session has been established and the Link-State/SPF
adddress family capability has been exchanged [RFC4790] on the
corresponding session.  If the session goes down, the NLRI will be
withdrawn.

## 2.2.  BGP Peering Between Directly Connected Network Nodes

In this model, BGP speakers peer with all directly connected network
nodes but the sessions may be multi-hop and the direct connection
discovery and liveliness detection for those connections are
independent of the BGP protocol.  How this is accomplished is outside
the scope of this document.  Consequently, there will be a single
session even if there are multiple direct connections between BGP
speakers.  For the purposes of BGP SPF, Link NLRI is advertised as
long as a BGP session has been established, the Link-State/SPF

address family capability has been exchanged [RFC4790] and the
corresponding link is up and considered operational.

## 2.3.  BGP Peering in Route-Reflector or Controller Topology

In this model, BGP speakers peer solely with one or more Route
Reflectors [RFC4456] or controllers.  As in the previous model,
direct connection discovery and liveliness detection for those
connections are done outside the BGP protocol.  For the purposes of
BGP SPF, Link NLRI is advertised as long as the corresponding link is
up and considered operational.

## 3.  BGP-LS Shortest Path Routing (SPF) SAFI

In order to replace the Phase 1 and 2 decision functions of the
existing Decision Process with an SPF-based Decision Process and
streamline the Phase 3 decision functions in a backward compatible
manner, this draft introduces a couple AFI/SAFIs for BGP LS SPF
operation.  The BGP-LS-SPF (AF 16388 / SAFI TBD1) and BGP-LS-SPF-VPN
(AFI 16388 / SAFI TBD2) [RFC4790] are allocated by IANA as specified
in the Section 6.

## 4.  Extensions to BGP-LS

[RFC7752] describes a mechanism by which link-state and TE
information can be collected from networks and shared with external
components using BGP protocol.  It contains two parts: definition of
a new BGP NLRI that describes links, nodes, and prefixes comprising
IGP link-state information and definition of a new BGP path attribute
(BGP-LS attribute) that carries link, node, and prefix properties and
attributes, such as the link and prefix metric or auxiliary Router-
IDs of nodes, etc.

The BGP protocol will be used in the Protocol-ID field specified in
table 1 of [I-D.ietf-idr-bgpls-segment-routing-epe].  The local and
remote node descriptors for all NLRI will be the BGP Router-ID (TLV
516) and either the AS Number (TLV 512) [RFC7752] or the BGP
Confederation Member (TLV 517)
[I-D.ietf-idr-bgpls-segment-routing-epe].  However, if the BGP
Router-ID is known to be unique within the BGP Routing domain, it can
be used as the sole descriptor.

## 4.1.  Node NLRI Usage and Modifications

The SPF capability is a new Node Attribute TLV that will be added to
those defined in table 7 of [RFC7752].  The new attribute TLV will
only be applicable when BGP is specified in the Node NLRI Protocol ID

field.  The TBD TLV type will be defined by IANA.  The new Node
Attribute TLV will contain a single octet SPF algorithm field:

```
    0                   1                   2                   3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |             Type              |             Length            |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   | SPF Algorithm |
   +-+-+-+-+-+-+-+-+
```

 The SPF Algorithm may take the following values:

   1 - Normal SPF
   2 - Strict SPF


When computing the SPF for a given BGP routing domain, only BGP nodes
advertising the SPF capability attribute will be included the
Shortest Path Tree (SPT).

## 4.2.  Link NLRI Usage

The criteria for advertisement of Link NLRI are discussed in
Section 2.

Link NLRI is advertised with local and remote node descriptors as
described above and unique link identifiers dependent on the
addressing.  For IPv4 links, the links local IPv4 (TLV 259) and
remote IPv4 (TLV 260) addresses will be used.  For IPv6 links, the
local IPv6 (TLV 261) and remote IPv6 (TLV 262) addresses will be
used.  For unnumbered links, the link local/remote identifiers (TLV
258) will be used.  For links supporting having both IPv4 and IPv6
addresses, both sets of descriptors may be included in the same Link
NLRI.  The link identifiers are described in table 5 of [RFC7752].

The link IGP metric attribute TLV (TLV 1095) as well as any others
required for non-SPF purposes SHOULD be advertised.  Algorithms such
as setting the metric inversely to the link speed as done in the OSPF
MIB [RFC4750] may be supported.  However, this is beyond the scope of
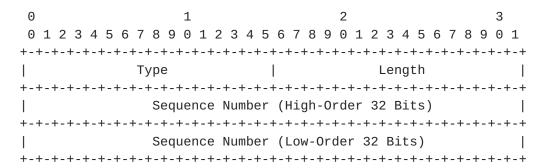this document.

## 4.3.  Prefix NLRI Usage

Prefix NLRI is advertised with a local descriptor as described above
and the prefix and length used as the descriptors (TLV 265) as
described in [RFC7752].  The prefix metric attribute TLV (TLV 1155)
as well as any others required for non-SPF purposes SHOULD be

advertised.  For loopback prefixes, the metric should be 0.  For non-loopback, the setting of the metric is beyond the scope of this document.

### 4.4.  BGP-LS Attribute Sequence-Number TLV

A new BGP-LS Attribute TLV to BGP-LS NLRI types is defined to assure the most recent version of a given NLRI is used in the SPF computation.  The TBD TLV type will be defined by IANA.  The new BGP-LS Attribute TLV will contain an 8 octet sequence number.  The usage of the Sequence Number TLV is described in Section 5.1.

```
  0                   1                   2                   3
  0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |             Type              |            Length             |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |             Sequence Number (High-Order 32 Bits)             |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |             Sequence Number (Low-Order 32 Bits)              |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

Sequence Number

The 64-bit strictly increasing sequence number is incremented for every version of BGP-LS NLRI originated.  BGP speakers implementing this specification MUST use available mechanisms to preserve the sequence number's strictly increasing property for the deployed life of the BGP speaker (including cold restarts).  One mechanism for accomplishing this would be to use the high-order 32 bits of the sequence number as a wrap/boot count that is incremented anytime the BGP Router router loses its sequence number state or the low-order 32 bits wrap.

When incrementing the sequence number for each self-originated NLRI, the sequence number should be treated as an unsigned 64-bit value. If the lower-order 32-bit value wraps, the higher-order 32-bit value should be incremented and saved in non-volatile storage.  If by some chance the BGP Speaker is deployed long enough that there is a possibility that the 64-bit sequence number may wrap or a BGP Speaker completely loses its sequence number state (e.g, the BGP speaker hardware is replaced), the phase 1 decision function (see Section 5.1) rules should insure convergance, albeit, not immediately.

## 5.  Decision Process with SPF Algorithm

The Decision Process described in [RFC4271] takes place in three
distinct phases.  The Phase 1 decision function of the Decision
Process is responsible for calculating the degree of preference for
each route received from a Speaker's peer.  The Phase 2 decision
function is invoked on completion of the Phase 1 decision function
and is responsible for choosing the best route out of all those
available for each distinct destination, and for installing each
chosen route into the Loc-RIB.  The combination of the Phase 1 and 2
decision functions is also known as a Path vector algorithm.

When BGP-LS-SPF NLRI is received, all that is required is to
determine whether it is the best-path by examining the Node-ID and
sequence number as described in Section 5.1.  If the best-path NLRI
had changed, it will be advertised to other BGP-LS-SPF peers.  If the
attributes have changed (other than the sequence number), a BGP SPF
calculation will be scheduled.  However, a changed best-path can be
advertised to other peer immediately and propagation of changes can
approach IGP convergence times.

The SPF based Decision process starts with selecting only those Node
NLRI whose SPF capability TLV matches with the local BGP speaker's
SPF capability TLV value.  Since Link-State NLRI always contains the
local descriptor [RFC7752], it will only be originated by a single
BGP speaker in the BGP routing domain.  These selected Node NLRI and
their Link/Prefix NLRI are used to build a directed graph during the
SPF computation.  The best paths for BGP prefixes are installed as a
result of the SPF process.

The Phase 3 decision function of the Decision Process [RFC4271] is
also simplified since under normal SPF operation, a BGP speaker would
advertise the NLRI selected for the SPF to all BGP peers with the
BGP-LS/BGP-SPF AFI/SAFI.  Application of policy would not be
prevented but would normally not be necessary.

## 5.1.  Phase-1 BGP NLRI Selection

The rules for NLRI selection are greatly simplified from [RFC4271].

1.  If the NLRI is received from the BGP speaker originating the NLRI
    (as determined by the comparing BGP Router ID in the NLRI Node
    identifiers with the BGP speaker Router ID), then it is preferred
    over the same NLRI from non-originators.

2.  If the Sequence-Number TLV is present in the BGP-LS Attribute,
    then the NLIR with the most recent, i.e., highest sequence number
    is selected.  BGP-LS NLRI with a Sequence-Number TLV will be

considered more recent than NLRI without a BGP-LS or a BGP-LS
Attribute that doesn't include the Sequence-Number TLV.

3.  The final tie-breaker is the NLRI from the BGP Speaker with the
    numerically largest BGP Router ID.

The modified Decision Process with SPF algorithm uses the metric from
Link and Prefix NLRI Attribute TLVs [RFC7752].  As a result, any
attributes that would influence the Decision process defined in
[RFC4271] like ORIGIN, MULTI_EXIT_DISC, and LOCAL_PREF attributes are
ignored by the SPF algorithm.  Furthermore, the NEXT_HOP attribute
value is preserved and validated but otherwise ignored during the SPF
or best-path.

## 5.2.  Dual Stack Support

The SPF based decision process operates on Node, Link, and Prefix
NLRIs that support both IPv4 and IPv6 addresses.  Whether to run a
single SPF instance or multiple SPF instances for separate AFs is a
matter of a local implementation.  Normally, IPv4 next-hops are
calculated for IPv4 prefixes and IPv6 next-hops are calculated for
IPv6 prefixes.  However, an interesting use-case is deployment of
[RFC5549] where IPv6 link-local next-hops are calculated for both
IPv4 and IPv6 prefixes.  As stated in Section 1, support for Multiple
Topology Routing (MTR) is an area for future study.

## 5.3.  NEXT_HOP Manipulation

A BGP speaker that supports SPF extensions MAY interact with peers
that don't support SPF extensions.  If the BGP Link-State address
family is advertised to a peer not supporting the SPF extensions
described herein, then the BGP speaker MUST conform to the NEXT_HOP
rules mentioned in [RFC4271] when announcing the Link-State address
family routes to those peers.

All BGP peers that support SPF extensions would locally compute the
NEXT_HOP values as result of the SPF process.  As a result, the
NEXT_HOP attribute is always ignored on receipt.  However BGP
speakers should set the NEXT_HOP address according to the NEXT_HOP
attribute rules mentioned in [RFC4271].

## 5.4.  IPv4/IPv6 Unicast Address Family Interaction

While the BGP-LS SPF address family and the IPv4/IPv6 unicast address
families install routes into the same device routing tables, they
will operate independently much the same as OSPF and IS-IS would
operate today (i.e., "Ships-in-the-Night" mode).  There will be no
implicit route redistribution between the BGP address families.

However, implementation specific redistribution mechanisms SHOULD be
made available with the restriction that redistribution of BGP-LS SPF
routes into the IPv4 address family applies only to IPv4 routes and
redistribution of BGP-LS SPF route into the IPv6 address family
applies only to IPv6 routes.

Given the fact that SPF algorithms are based on the assumption that
all routers in the routing domain calculate the precisely the same
SPF tree and install the same set of routers, it is RECOMMENDED that
BGP-LS SPF IPv4/IPv6 routes be given priority by default when
installed into their respective RIBs.  In common implementations the
prioritization is governed by route preference or administrative
distance with lower being more preferred.

## 5.5.  NLRI Advertisement and Convergence

A local failure will prevent a link from being used in the SPF
calculation due to the IGP bi-directional connectivity requirment.
Consequently, local link failues should always be given priority over
updates (e.g., withdrawing all routes learned on a session) in order
to ensure the highest priority progation and optimal convergence.

Delaying the withdrawal of non-local routes is an area for further
study as more IGP-like mechanisms would be required to prevent usage
of stale NLRI.

## 5.6.  Error Handling

When a BGP speaker receives a BGP Update containing a malformed SPF
Capability TLV in the Node NLRI BGP-LS Attribute [RFC7752], it MUST
ignore the received TLV and the Node NLRI and not pass it to other
BGP peers as specified in [RFC7606].  When discarding a Node NLRI
with malformed TLV, a BGP speaker SHOULD log an error for further
analysis.

## 6.  IANA Considerations

This document defines a couple AFI/SAFIs for BGP LS SPF operation and
requests IANA to assign the BGP-LS-SPF AFI 16388 / SAFI TBD1 and the
BGP-LS-SPF-VPN AFI 16388 / SAFI TBD2 as described in [RFC4750].

This document also defines two attribute TLV for BGP LS NLRI.  We
request IANA to assign TLVs for the SPF capability and the Sequence
Number from the "BGP-LS Node Descriptor, Link Descriptor, Prefix
Descriptor, and Attribute TLVs" Registry.  Additionally, IANA is
requested to create a new registry for "BGP-LS SPF Capability
Algorithms" for the value of the algorithm both in the BGP-LS Node

Attribute TLV and the BGP SPF Capability.  The initial assignments are:

```
+-------------+----------------------------------+
| Value(s)    | Assignment Policy                |
+-------------+----------------------------------+
| 0           | Reserved (not to be assigned)    |
|             |                                  |
| 1           | SPF                              |
|             |                                  |
| 2           | Strict SPF                       |
|             |                                  |
| 3-254       | Unassigned (IETF Review)         |
|             |                                  |
| 255         | Reserved (not to be assigned)    |
+------------+----------------------------------+
```

                    BGP-LS SPF Capability Algorithms

## 7.  Security Considerations

This extension to BGP does not change the underlying security issues inherent in the existing [RFC4724] and [RFC4271].

## 7.1.  Acknowledgements

The authors would like to thank .... for the review and comments.

## 7.2.  Contributorss

In addition to the authors listed on the front page, the following co-authors have contributed to the document.

  Derek Yeung
  Arrcus, Inc.
  derek@arrcus.com

  Gunter Van De Velde
  Nokia
  gunter.van_de_velde@nokia.com

  Abhay Roy
  Cisco Systems
  akr@cisco.com

  Venu Venugopal
  Cisco Systems
  venuv@cisco.com

## 8. References

### 8.1. Normative References

[I-D.ietf-idr-bgpls-segment-routing-epe]
          Previdi, S., Filsfils, C., Patel, K., Ray, S., and J.
          Dong, "BGP-LS extensions for Segment Routing BGP Egress
          Peer Engineering", draft-ietf-idr-bgpls-segment-routing-
          epe-14 (work in progress), December 2017.

[RFC2119]  Bradner, S., "Key words for use in RFCs to Indicate
          Requirement Levels", BCP 14, RFC 2119,
          DOI 10.17487/RFC2119, March 1997, <https://www.rfc-
          editor.org/info/rfc2119>.

[RFC4271]  Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A
          Border Gateway Protocol 4 (BGP-4)", RFC 4271,
          DOI 10.17487/RFC4271, January 2006, <https://www.rfc-
          editor.org/info/rfc4271>.

[RFC7606]  Chen, E., Ed., Scudder, J., Ed., Mohapatra, P., and K.
          Patel, "Revised Error Handling for BGP UPDATE Messages",
          RFC 7606, DOI 10.17487/RFC7606, August 2015,
          <https://www.rfc-editor.org/info/rfc7606>.

[RFC7752]  Gredler, H., Ed., Medved, J., Previdi, S., Farrel, A., and
          S. Ray, "North-Bound Distribution of Link-State and
          Traffic Engineering (TE) Information Using BGP", RFC 7752,
          DOI 10.17487/RFC7752, March 2016, <https://www.rfc-
          editor.org/info/rfc7752>.

[RFC7938]  Lapukhov, P., Premji, A., and J. Mitchell, Ed., "Use of
          BGP for Routing in Large-Scale Data Centers", RFC 7938,
          DOI 10.17487/RFC7938, August 2016, <https://www.rfc-
          editor.org/info/rfc7938>.

### 8.2. Information References

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

[RFC4456]  Bates, T., Chen, E., and R. Chandra, "BGP Route
          Reflection: An Alternative to Full Mesh Internal BGP
          (IBGP)", RFC 4456, DOI 10.17487/RFC4456, April 2006,
          <https://www.rfc-editor.org/info/rfc4456>.

   [RFC4724]  Sangli, S., Chen, E., Fernando, R., Scudder, J., and Y.
              Rekhter, "Graceful Restart Mechanism for BGP", RFC 4724,
              DOI 10.17487/RFC4724, January 2007, <https://www.rfc-
              editor.org/info/rfc4724>.

   [RFC4750]  Joyal, D., Ed., Galecki, P., Ed., Giacalone, S., Ed.,
              Coltun, R., and F. Baker, "OSPF Version 2 Management
              Information Base", RFC 4750, DOI 10.17487/RFC4750,
              December 2006, <https://www.rfc-editor.org/info/rfc4750>.

   [RFC4790]  Newman, C., Duerst, M., and A. Gulbrandsen, "Internet
              Application Protocol Collation Registry", RFC 4790,
              DOI 10.17487/RFC4790, March 2007, <https://www.rfc-
              editor.org/info/rfc4790>.

   [RFC4915]  Psenak, P., Mirtorabi, S., Roy, A., Nguyen, L., and P.
              Pillay-Esnault, "Multi-Topology (MT) Routing in OSPF",
              RFC 4915, DOI 10.17487/RFC4915, June 2007,
              <https://www.rfc-editor.org/info/rfc4915>.

   [RFC5286]  Atlas, A., Ed. and A. Zinin, Ed., "Basic Specification for
              IP Fast Reroute: Loop-Free Alternates", RFC 5286,
              DOI 10.17487/RFC5286, September 2008, <https://www.rfc-
              editor.org/info/rfc5286>.

   [RFC5549]  Le Faucheur, F. and E. Rosen, "Advertising IPv4 Network
              Layer Reachability Information with an IPv6 Next Hop",
              RFC 5549, DOI 10.17487/RFC5549, May 2009,
              <https://www.rfc-editor.org/info/rfc5549>.

Authors' Addresses

   Keyur Patel
   Arrcus, Inc.

   Email: keyur@arrcus.com


   Acee Lindem
   Cisco Systems
   301 Midenhall Way
   Cary, NC  27513
   USA

   Email: acee@cisco.com

Shawn Zandi
Linkedin
222 2nd Street
San Francisco, CA  94105
USA

Email: szandi@linkedin.com


Wim Henderickx
Nokia
Antwerp
Belgium

Email: wim.henderickx@nokia.com