

Routing Area
Internet-Draft
Intended status: Informational
Expires: 19 August 2026

A. Bashandy, Ed.
HPE
C. Filsfils
Cisco Systems
P. Mohapatra
Sproute Networks
Y. Qu, Ed.
Futurewei Technologies
15 February 2026

BGP Prefix Independent Convergence
draft-ietf-rtgwg-bgp-pic-23

Abstract

In a network comprising thousands of BGP peers exchanging millions of routes, it is desirable to restore traffic after failure in a time period that does not depend on the number of BGP prefixes.

This document describes an architecture by which traffic can be re-routed to Equal Cost Multi-Path (ECMP) or pre-calculated backup paths in a timeframe that does not depend on the number of BGP prefixes. The objective is achieved through organizing the forwarding data structures in a hierarchical manner and sharing forwarding elements among the maximum possible number of routes. The described technique yields prefix independent convergence while ensuring incremental deployment, complete automation, and zero management and provisioning effort. It is noteworthy to mention that the benefits of BGP Prefix Independent Convergence (BGP-PIC) are hinged on the existence of more than one path whether as ECMP or primary-backup.

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 <https://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 19 August 2026.

Copyright Notice

Copyright (c) 2026 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	3
1.1. Terminology	4
2. Overview	5
2.1. Dependency	6
2.1.1. Hierarchical Hardware FIB (Forwarding Information Base)	6
2.1.2. Availability of Precomputed Backup Paths	6
2.2. BGP-PIC Illustration	7
3. Constructing the Shared Hierarchical Forwarding Chain	9
3.1. Constructing the BGP-PIC Forwarding Chain	10
3.2. Example: Primary-Backup path Scenario	10
4. Forwarding Behavior	11
5. Handling Platforms with Limited Levels of Hierarchy	13
6. Forwarding Chain Adjustment at a Failure	13
6.1. BGP-PIC core	13
6.2. BGP-PIC edge	14
6.2.1. Adjusting Forwarding Chain in egress node failure	14
6.2.2. Adjusting Forwarding Chain on PE-CE link Failure	15
6.3. Handling Failures for Flattened Forwarding Chains	16
7. Operational Properties	17
7.1. Failure Coverage	17
7.2. Convergence Characteristics	18
7.3. Fast Local Repair	18
7.4. Configuration Free	18
7.5. Incremental Deployment	18
8. Security Considerations	19
9. IANA Considerations	19
10. References	19
10.1. Normative References	19
10.2. Informative References	19
Appendix A. Acknowledgments	21

Appendix B. Handling Platforms with Limited Levels of Hierarchy	21
Appendix C. Example: Flattening a forwarding chain.	23
Appendix D. Perspective	30
Authors' Addresses	31

1. Introduction

BGP speakers exchange reachability information about prefixes [RFC4271]. For labeled address families, an edge router assigns local labels to prefixes and associates the local label with each advertised prefix using technologies such as L3VPN [RFC4364], 6PE [RFC4798], and Softwire [RFC5565] using BGP label unicast (BGP-LU) technique [RFC8277]. A BGP speaker then applies the path selection steps to choose the best route. In modern networks, it is not uncommon to have a prefix reachable via multiple edge routers. Multiple techniques have been described to allow for BGP to advertise more than one path for a given prefix [I-D.ietf-idr-best-external][RFC7911][RFC6774], whether in the form of equal cost multipath or primary-backup. Another common and widely deployed scenario is L3VPN with multi-homed VPN sites with unique Route Distinguisher.

This document describes a hierarchical and shared forwarding chain organization that allows traffic to be restored to a pre-calculated alternative equal cost path or backup path in a time period that does not depend on the number of BGP prefixes. The technique relies on internal router behavior that is completely transparent to the operator and can be incrementally deployed and enabled with zero operator intervention. In other words, once it is implemented and deployed on a router, nothing is required from the operator to make it work. It is noteworthy to mention that this document describes a Forwarding Information Base (FIB) architecture that can be implemented in both hardware and/or software, although we refer to hardware implementation in most of the cases because of the additional complexity and performance requirements associated with hardware implementations.

It should be noted that although BGP is used for routes calculation in this document, the underlying principles of hierarchical forwarding, recursive resolution are not BGP specific. These mechanisms apply equally to routes computed by other routing protocols as well. The benefits of BGP-PIC are tied to the forwarding plane design rather than to the BGP protocol.

1.1. Terminology

This section defines the terms used in this document.

- * BGP-LU: BGP Label Unicast. Refers to using BGP to advertise the binding of an address prefix to one or more MPLS labels as in [RFC8277].
- * BGP prefix: A set of destination as an IP prefix with route learned through BGP as described in [RFC4271].
- * IGP prefix: A prefix that is learned via an Interior Gateway Protocol (IGP), such as OSPF and IS-IS.
- * ePE: Egress PE [RFC4364].
- * iPE: Ingress PE [RFC4364].
- * Path: One specific candidate way to reach the destination in a route [RFC4271]. It's a sequence of nodes or links from the source to the destination. The nodes may not be directly connected.
- * Recursive path: The next-hop of a path is an IP without the outgoing interface. it requires the router to look up the next-hop IP in the routing table (recursion) until it finds a directly connected or attached next-hop.
- * Non-recursive path: A path consisting of the IP address of a directly connected next-hop and outgoing interface.
- * Adjacency: The layer 2 encapsulation leading to the layer 3 directly connected next-hop. An adjacency is identified by a next-hop and an outgoing interface
- * Primary path: A recursive or non-recursive path that can be used for forwarding. A prefix can have more than one primary path.
- * Backup path: A recursive or non-recursive path that can be used only after some or all primary paths become unreachable.
- * Leaf: A container data structure for a prefix or local label. Alternatively, it is the data structure that contains prefix specific information.
- * IP leaf: The leaf corresponding to an IPv4 or IPv6 prefix.

- * Label leaf. The leaf corresponding to a locally allocated label such as the VPN label on an egress PE [RFC4364].
- * Pathlist: An array of paths used by one or more prefixes to forward traffic to destination(s) covered by an IP prefix. Each path in the pathlist carries its "path-index" that identifies its position in the array of paths. In general the value of the path-index in a path is the same as its position in the pathlist, except in the case outlined in Section 5. For example the 3rd path may carry a path-index value of 1. A pathlist may contain a mix of primary and backup paths.
- * OutLabel-List: Each labeled prefix is associated with an OutLabel-List. The OutLabel-List is an array of one or more outgoing labels and/or label actions where each label or label action has 1-to-1 correspondence to a path in the pathlist. Label actions are: push (add) the label as specified in [RFC3031], pop (remove) the label as specified in [RFC3031], swap (replace) the incoming label with the label in the OutLabel-List entry, or don't push anything at all in case of "unlabeled". The prefix may be an IGP or BGP prefix.
- * Forwarding chain: It is a compound data structure consisting of multiple connected blocks that a forwarding engine walks one block at a time to forward the packet out of an interface. Section 2.2 explains an example of a forwarding chain. Subsequent sections provide additional examples
- * Dependency: An object X is said to be a dependent or child of object Y if there is at least one forwarding chain where the forwarding engine must visit the object X before visiting the object Y in order to forward a packet. Note that if object X is a child of object Y, then Y cannot be deleted unless object X is no longer a dependent/child of object Y.
- * ASN: Autonomous System Number.

2. Overview

The idea of BGP-PIC is based on the following two pillars to make convergence independent of the number of prefixes:

- * A shared hierarchical forwarding chain: Multiple prefixes reference common next-hop and path objects arranged in a hierarchy, so that changes to a single shared object affect all dependent prefixes simultaneously.

- * A forwarding plane with multiple levels of indirection: The forwarding plane supports recursive resolution and pointer-based forwarding entries, allowing failover by updating a small number of shared objects rather than per-prefix state.

A forwarding plane with shared, hierarchical forwarding chains with maximal object reuse can reroute a large number of destinations by modifying only a small set of shared objects. This enables convergence in a time frame that does not depend on the number of affected destinations. For example, if an IGP prefix used to resolve a recursive next-hop changes, there is no need to update the potentially large number of BGP NLRI's that reference that next-hop.

2.1. Dependency

This section describes the required functionalities in the forwarding and control planes to support BGP-PIC as described in this document.

2.1.1. Hierarchical Hardware FIB (Forwarding Information Base)

BGP-PIC requires forwarding hardware that supports a hierarchical FIB. When a packet's destination address matches a BGP prefix, the forwarding plane performs recursive lookups through successive levels of indirection until a resolving adjacency is reached. Section 4 provides further details on the packet forwarding process.

For platforms that support only a limited number of levels of indirection, a necessary trade-off approach is to flatten forwarding dependencies when programming BGP destinations into the hardware FIB. In this case, recursive resolution is resolved at programming time, potentially eliminating both BGP pathlist and IGP pathlist lookups during forwarding.

While flattening reduces the number of memory accesses per packet, it comes at the cost of increased hardware FIB memory usage as flattening reduces sharing and results in greater duplication of forwarding entries, reduced ECMP and BGP-PIC properties as fewer pathlists are available.

Appendix B describes the flattening approach in more detail for hardware platforms with a limited number of supported indirection levels.

2.1.2. Availability of Precomputed Backup Paths

BGP-PIC requires backup paths so that traffic can be immediately redirected in the forwarding plane when a next hop fails, without reprocessing individual BGP prefixes.

Backup paths are calculated before any failure and installed in the FIB along with the primary path. Because many prefixes share the same next hop, a failure only requires switching that next hop to its back.

The BGP distribution of multiple paths is available thanks to the following BGP mechanisms: Add-Path [RFC7911], BGP Best-External [I.D.ietf-idr-best-external], diverse path [RFC6774], and the frequent use in VPN deployments of different VPN RD's per PE. Another option to learn multiple BGP next-hops/paths is to receive IBGP paths from multiple BGP RRs [RFC9107] selecting a different path as best. It is noteworthy to mention that the availability of another BGP path does not mean that all failure scenarios can be covered by simply forwarding traffic to the available secondary path. The discussion of how to cover various failure scenarios is beyond the scope of this document.

2.2. BGP-PIC Illustration

To illustrate the two pillars above as well as the platform dependency, this document will use an example of a multihomed L3VPN prefix in a BGP-free core running LDP [RFC5036] or segment routing over MPLS forwarding plane [RFC8660].

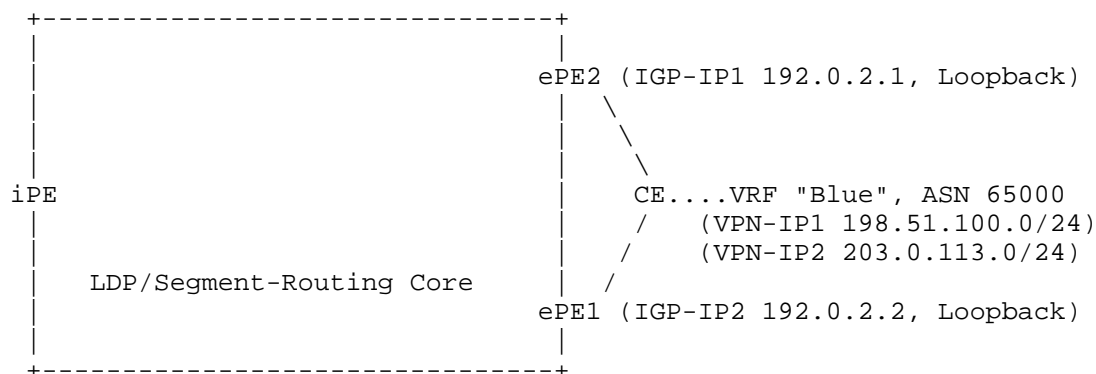


Figure 1: VPN prefix reachable via multiple PEs

Referring to Figure 1, suppose the iPE (the ingress PE) receives NLRIs for the VPN prefixes VPN-IP1 and VPN-IP2 from two egress PEs, ePE1 and ePE2 with next-hop BGP-NH1 (192.0.2.1) and BGP-NH2 (192.0.2.2), respectively. Assume that ePE1 advertise the VPN labels VPN-L11 and VPN-L12 while ePE2 advertise the VPN labels VPN-L21 and VPN-L22 for VPN-IP1 and VPN-IP2, respectively. Suppose that BGP-NH1 and BGP-NH2 are resolved via the IGP prefixes IGP-IP1 and IGP-IP2, where each happen to have 2 equal cost paths with IGP-NH1 and IGP-NH2 reachable via the interfaces I1 and I2 on iPE, respectively.

pathlist needs to be modified. Likewise, due to the hierarchical structure of the forwarding chain, it is possible to make modifications to the IGP routes without having to make any changes to the BGP NLRI's. For example, if the interface "I2" goes down, only the shared IGP pathlist needs to be updated, but none of the IGP prefixes sharing the IGP pathlist nor the BGP NLRI's using the IGP prefixes for resolution need to be modified.

Figure 2 can also be used to illustrate the second BGP-PIC pillar. Having a deep forwarding chain such as the one illustrated in Figure 2 requires a forwarding plane that is capable of accessing multiple levels of indirection in order to calculate the outgoing interface(s) and next-hops(s). While a deeper forwarding chain minimizes the re-convergence time on topology change, there will always exist platforms with limited capabilities and hence imposing a limit on the depth of the forwarding chain. Section 5 describes how to gracefully trade off convergence speed with the number of hierarchical levels to support platforms with different capabilities.

Another example using IPv6 addresses can be something like the following:

```
65000: 2001:DB8:1::/48
    via ePE1 (65000: 2001:DB8:192::1), VPN Label: VPN6-L11
    via ePE2 (65000: 2001:DB8:192::2), VPN Label: VPN6-L21

65000: 2001:DB8:2:/48
    via ePE1 (65000: 2001:DB8:192::1), VPN Label: VPN6-L12
    via ePE2 (65000: 2001:DB8:192::2), VPN Label: VPN6-L22

65000: 2001:DB8:192::1/128
    via Core, Label:      IGP6-L11
    via Core, Label:      IGP6-L12

65000: 2001:DB8:192::2/128
    via Core, Label:      IGP6-L21
    via Core, Label:      IGP6-L22
```

The same hierarchical forwarding chain described can be constructed for IPv6 addresses/prefixes.

3. Constructing the Shared Hierarchical Forwarding Chain

This section describes how the forwarding chain is constructed using a hierarchical shared model, as introduced in Section 2. Section 3.1 details the construction steps, and Section 3.2 provides an illustrative example.

3.1. Constructing the BGP-PIC Forwarding Chain

The forwarding chain is built using the following steps:

- (1) Prefix arrival in FIB. The prefix contains one or more outgoing paths. For certain labeled prefixes, such as L3VPN [RFC4364] prefixes, each path may be associated with an outgoing label and the prefix itself may be assigned a local label. The list of outgoing paths defines a pathlist.
- (2) Pathlist lookup/creation. If such pathlist does not already, then the FIB manager (software or hardware entity responsible for managing the FIB) creates a new pathlist, otherwise the existing pathlist with the same list of paths exist (the pathlist may already exist because there is another pic-route that is already using the same list of paths) is used.
- (3) Register prefix dependency. The BGP prefix is added as a dependent of the pathlist.
- (4) Resolve pathlist entries. The forwarding chain is completed by resolving the paths of the pathlist. A BGP path usually consists of a next-hop. The next-hop is resolved by finding a matching prefix reachable via IGP or other protocols.

The end result is a hierarchical shared forwarding chain where the BGP pathlist is shared by all BGP prefixes that use the same list of paths and the IGP prefix is shared by all pathlists that have a path resolving via that IGP prefix.

3.2. Example: Primary-Backup path Scenario

Consider the egress PE ePE1 in the case of the multi-homed VPN prefixes shown in Figure 1. Suppose ePE1 determines that the primary path is the external path, while the backup path is the IBGP path to the other PE ePE2 with next-hop BGP-NH2. ePE1 constructs the forwarding chain depicted in Figure 3. The figure shows only a single VPN prefix for simplicity. But all prefixes that are multihomed to ePE1 and ePE2 share the BGP pathlist.

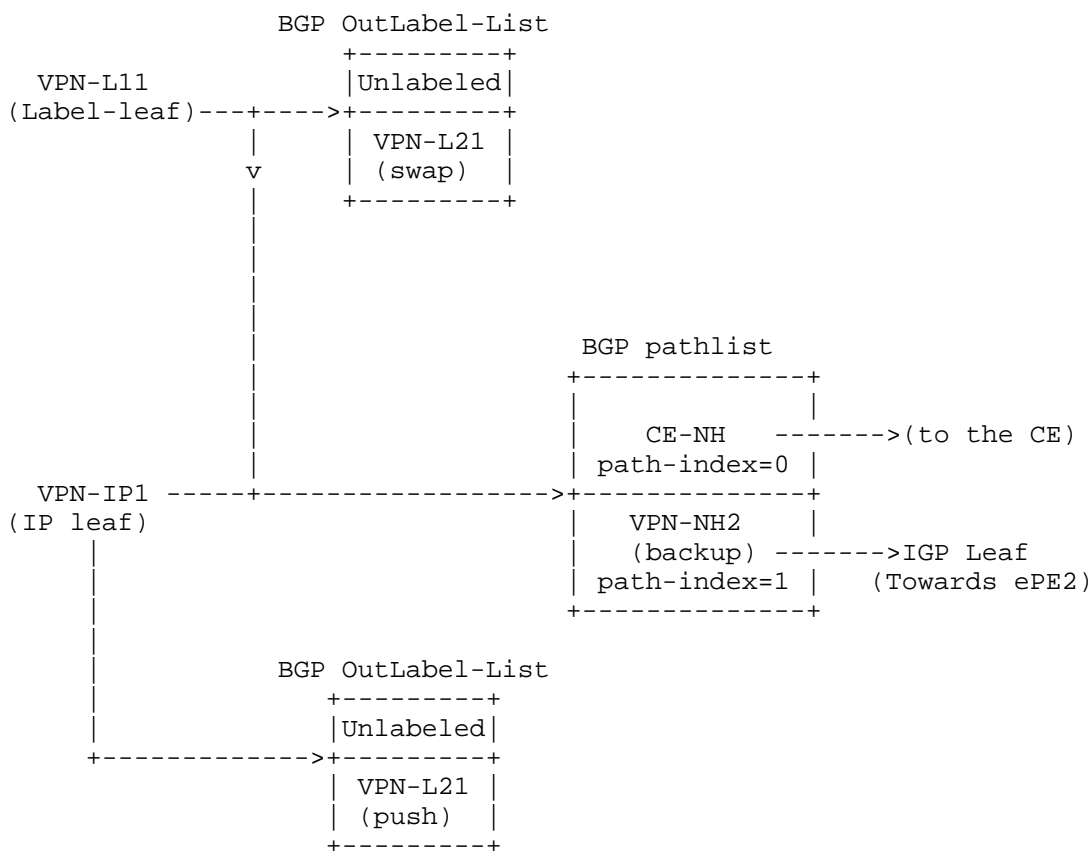


Figure 2: VPN Prefix Forwarding Chain with eiBGP paths on egress PE

The example depicted in Figure 3 differs from the example in Figure 2 in two main aspects. First, as long as the primary path towards the CE (external path) can be used for forwarding, it will be the only path used for forwarding while the OutLabel-List contains both the unlabeled (primary path) and the VPN label (backup path) advertised by the backup path ePE2. The second aspect is presence of the label leaf corresponding to the VPN prefix. This label leaf is used to match VPN traffic arriving from the core. Note that the label leaf shares the pathlist with the IP prefix.

4. Forwarding Behavior

This section explains how the forwarding plane uses the hierarchical shared forwarding chain to forward a packet.

When a packet arrives at a router, assume it matches a leaf. If not, the packet is handled according to the local policy (such as silently dropping the packet), which is beyond the scope of this document. A labeled packet matches a label leaf while an IP packet matches an IP leaf. The forwarding engine walks the forwarding chain starting from the leaf until the walk terminates on an adjacency. Thus when a packet arrives, the chain is walked as follows:

1. Lookup the leaf based on the destination address or the label at the top of the packet.
2. Retrieve the parent pathlist of the leaf.
3. Pick an outgoing path "Pi" from the list of resolved pic- paths in the pathlist. The method by which the outgoing path is picked is beyond the scope of this document (e.g. flow- preserving hash exploiting entropy within the MPLS stack and IP header). Let the "path-index" of the outgoing path "Pi" be "j". Remember that, as described in the definition of the term pathlist in Section 1.1, the path-index of a path may not always be identical the position of the path in the pathlist.
4. If the prefix is labeled, use the "path-index" "j" to retrieve the label "Lj" stored position j in the OutLabel-List and apply the label action of the label on the packet (e.g. for VPN label on the ingress PE, the label action is "push"). As mentioned in Section 1.1 the value of the "path-index" stored in the pic- path may not necessarily be the same value of the location of the path in the pathlist.
5. If the chosen path "Pi" is recursive, move to its parent prefix and go to step 2.
6. If the chosen path is non-recursive move to its parent adjacency.
7. Encapsulate the packet in the layer string specified by the adjacency and send the packet out.

Let's apply the above forwarding steps to the forwarding chain depicted in Figure 2 in Section 2. Suppose a packet arrives at ingress PE iPE from an external neighbor. Assume the packet matches the VPN prefix VPN-IP1. While walking the forwarding chain, the forwarding engine applies a hashing algorithm to choose the path and the hashing at the BGP level chooses the first path in the BGP pathlist while the hashing at the IGP level yields the second path in the IGP pathlist. In that case, the packet will be sent out of interface I2 with the label stack "IGP-L12,VPN-L11".

5. Handling Platforms with Limited Levels of Hierarchy

This section describes the construction of the forwarding chain if a platform does not support the number of recursion levels required to resolve the NLRIs. There are two main design objectives.

- * Being able to reduce the number of hierarchical levels from any arbitrary value to a smaller arbitrary value that can be supported by the forwarding engine.
- * Minimal modifications to the forwarding algorithm due to such reduction.

Appendix B provides details on how to handle limited hardware capabilities.

6. Forwarding Chain Adjustment at a Failure

The hierarchical and shared structure of the forwarding chain explained in the previous section allows modifying a small number of forwarding chain objects to re-route traffic to a pre-calculated equal-cost or backup path without the need to modify the possibly very large number of BGP prefixes. This section goes over various core and edge failure scenarios to illustrate how the FIB manager can utilize the forwarding chain structure to achieve BGP prefix independent convergence.

6.1. BGP-PIC core

This section describes the adjustments to the forwarding chain when a core link or node fails but the BGP next-hop remains reachable.

There are two case: remote link failure and attached link failure. Node failures are treated as link failures.

When a remote link or node fails, the IGP on the ingress PE receives an advertisement indicating a topology change so IGP re-converges to either find a new next-hop and/or outgoing interface or remove the path completely from the IGP prefix used to resolve BGP next-hops. IGP and/or LDP download the modified IGP leaves with modified outgoing labels for the labeled core.

When a local link fails, FIB manager detects the failure almost immediately. The FIB manager marks the impacted path(s) as unusable so that only useable paths are used to forward packets. Hence only IGP pathlists with paths using the failed local link need to be modified. All other pathlists are not impacted. Note that in this

particular case there is no need to backwalk (walk back the forwarding chain) to IGP leaves to adjust the OutLabel-Lists because FIB can rely on the path-index stored in the useable paths in the pathlist to pick the right label.

It is noteworthy to mention that because FIB manager modifies the forwarding chain starting from the IGP leaves only. BGP pathlists and leaves are not modified. Hence traffic restoration occurs within the time frame of IGP convergence, and, for local link failure, assuming a backup path has been precomputed, within the timeframe of local detection (e.g. 50ms). Examples of solutions that can precompute backup paths are IP FRR [RFC5714] remote LFA [RFC7490], TI-LFA [I-D.ietf-rtgwg-segment-routing-ti-lfa] and MRT [RFC7812] or EBGW path having a backup path [bonaventure].

Let's apply the procedure mentioned in this subsection to the forwarding chain depicted in Figure 2. Suppose a remote link failure occurs and impacts the first ECMP IGP path to the remote BGP next-hop. Upon IGP convergence, the IGP pathlist used by the BGP next-hop is updated to reflect the new topology (one path instead of two) and the new forwarding state is immediately available to all dependent BGP prefixes. The same behavior would occur if the failure was local such as an interface going down. As soon as the IGP convergence is complete for the BGP next-hop IGP pic-route, all its BGP depending routes benefit from the new pic-path. In fact, upon local failure, if LFA protection is enabled for the IGP route to the BGP next-hop and a backup path was pre-computed and installed in the pathlist, upon the local interface failure, the LFA backup path is immediately activated (e.g. sub- 50msec) and thus protection benefits all the depending BGP traffic through the hierarchical forwarding dependency between the routes.

6.2. BGP-PIC edge

This section describes the adjustments to the forwarding chains as a result of edge node or edge link failure.

6.2.1. Adjusting Forwarding Chain in egress node failure

When a node fails, IGP on neighboring core nodes send updates indicating that the edge node is no longer a direct neighbor. If the node that failed is an egress node, such as ePE1 and ePE2 in Figure 1, IGP running on an ingress node, such as iPE in Figure 1, converges and realizes that the egress node is no longer reachable. As such IGP on the ingress node instructs FIB to remove the IP and label leaves corresponding to the failed edge node from FIB. So FIB manager on the ingress node performs the following steps:

- * FIB manager deletes the IGP leaf corresponding to the failed edge node
- * FIB manager backwalks to all dependent BGP pathlists and marks that path using the deleted IGP leaf as unresolved
- * Note that there is no need to modify the possibly large number of BGP leaves because each path in the pathlist carries its pic- path index and hence the correct outgoing label will be picked. Consider for example the forwarding chain depicted in Figure 2. If the 1st BGP path becomes unresolved, then the forwarding engine will only use the second path for forwarding. Yet the path-index of that single resolved path will still be 1 and hence the label VPN-L21 will be pushed.

6.2.2. Adjusting Forwarding Chain on PE-CE link Failure

Suppose the link between an edge router and its external peer fails. There are two scenarios (1) the edge node attached to the failed link performs next-hop self (where BGP advertises the IP address of its own loopback as next-hop) and (2) the edge node attached to the failure advertises the IP address of the failed link as the next-hop attribute to its IBGP peers.

In the first case, the rest of IBGP peers will remain unaware of the link failure and will continue to forward traffic to the edge node until the edge node attached to the failed link withdraws the BGP prefixes. If the destination prefixes are multi-homed to another IBGP peer, say ePE2, then FIB manager on the edge router detecting the link failure applies the following steps to the forwarding chain (see Figure 3):

- * FIB manager backwalks to the BGP pathlists marks the path through the failed link to the external peer as unresolved.
- * Hence traffic will be forwarded using the backup path towards ePE2.
- * Labeled traffic arriving at the egress PE ePE1 matches the BGP label leaf.
 - The OutLabel-List attached to the BGP label leaf already contains an entry corresponding to the backup path.
 - The label entry in OutLabel-List corresponding to the internal path to backup egress PE has a swap action to the label advertised by the backup egress PE.

- For an arriving label packet (e.g. VPN), the top label is swapped with the label advertised by backup egress PE and the packet is sent towards that the backup egress PE.
- * Unlabeled traffic arriving at the egress PE ePE1 matches the BGP IP leaf
- The OutLabel-List attached to the BGP label leaf already contains an entry corresponding to the backup path.
 - The label entry in OutLabel-List corresponding to the internal path to backup egress PE has a push (instead of the swap action in for the labeled traffic case) action to the label advertised by the backup egress PE.
 - For an arriving IP packet, the label advertised by backup egress PE is pushed and the packet is sent towards that the backup egress PE.

In the second case where the edge router uses the IP address of the failed link as the BGP next-hop, the edge router will still perform the previous steps. But, unlike the case of next-hop self, the IGP on the failed edge node informs the rest of the IBGP peers that the IP address of the failed link is no longer reachable. Hence the FIB manager on IBGP peers will delete the IGP leaf corresponding to the IP prefix of the failed link. The behavior of the IBGP peers will be identical to the case of edge node failure outlined in Section 6.2.1.

It is noteworthy to mention that because the edge link failure is local to the edge router, sub-50 msec convergence can be achieved as described in [bonaventure].

Let's try to apply the case of next-hop self to the forwarding chain depicted in Figure 3. After failure of the link between ePE1 and CE, the forwarding engine will route traffic arriving from the core towards VPN-NH2 with path-index=1. A packet arriving from the core will contain the label VPN-L11 at top. The label VPN-L11 is swapped with the label VPN-L21 and the packet is forwarded towards ePE2.

6.3. Handling Failures for Flattened Forwarding Chains

As explained in the in Section 5 if the number of hierarchy levels of a platform cannot support the native number of hierarchy levels of a recursive forwarding chain, the instantiated forwarding chain is constructed by flattening two or more levels. Hence a 3-levels chain in Figure 5 is flattened into the 2-levels chain in Figure 6.

While reducing the benefits of BGP-PIC, flattening one hierarchy into a shallower hierarchy does not always result in a complete loss of the benefits of the BGP-PIC. To illustrate this fact suppose ASBR12 is no longer reachable in domain 1. If the platform supports the full hierarchy depth, the forwarding chain is the one depicted in Figure 5 and hence the FIB manager needs to backwalk one level to the pathlist shared by "ePE1" and "ePE2" and adjust it. If the platform supports 2 levels of hierarchy, then a useable forwarding chain is the one depicted in Figure 6. In that case, if ASBR12 is no longer reachable, the FIB manager has to backwalk to the two flattened pathlists and updates both of them.

The main observation is that the loss of convergence speed due to the loss of hierarchy depth depends on the structure of the forwarding chain itself. To illustrate this fact, let's take two extremes. Suppose the forwarding objects in level $i+1$ depend on the forwarding objects in level i . If every object on level $i+1$ depends on a separate object in level i , then flattening level i into level $i+1$ will not result in loss of convergence speed. Now let's take the other extreme. Suppose " n " objects in level $i+1$ depend on 1 object in level i . Now suppose FIB flattens level i into level $i+1$. If a topology change results in modifying the single object in level i , then FIB has to backwalk and modify " n " objects in the flattened level, thereby losing all the benefit of BGP-PIC. Experience shows that flattening forwarding chains usually results in moderate loss of BGP-PIC benefits. Further analysis is needed to corroborate and quantify this statement.

7. Operational Properties

7.1. Failure Coverage

BGP-PIC provides prefix-independent convergence for failures that affect shared forwarding dependencies, such as the loss of a next hop, an IGP path, or an adjacency used by multiple BGP prefixes. By precomputing and installing alternate forwarding paths and leveraging shared hierarchical forwarding objects, BGP-PIC enables traffic to be rerouted without requiring per-prefix BGP best-path recomputation.

Failures that do not impact shared forwarding objects, or that require BGP policy re-evaluation, may still rely on conventional BGP convergence behavior.

7.2. Convergence Characteristics

The primary convergence characteristic of BGP-PIC is that forwarding convergence time is independent of the number of affected BGP prefixes. Upon a failure, only a limited number of shared forwarding objects need to be updated. Compared with traditional BGP convergence, where forwarding updates scale with the number of impacted prefixes and may result in prolonged convergence in large routing tables.

7.3. Fast Local Repair

BGP-PIC enables forwarding repair that is independent of BGP control-plane convergence. Backup forwarding paths are computed and installed in advance, allowing the forwarding plane to redirect traffic immediately upon detection of a local failure.

When the failure is local (a local IGP next-hop failure or a local EBGp next-hop failure), a pre-computed and pre-installed backup is activated by a local-protection mechanism that does not depend on the number of BGP destinations impacted by the failure. Sub-50msec is thus possible even if millions of BGP prefixes are impacted.

When the failure is remote (a remote IGP failure not impacting the BGP next-hop or a remote BGP next-hop failure), an alternate pic-path is activated upon IGP convergence. All the impacted BGP destinations benefit from a working alternate path as soon as the IGP convergence occurs for their impacted BGP next-hop even if millions of BGP routes are impacted.

Appendix D puts the BGP-PIC benefits in perspective by providing some results using actual numbers.

7.4. Configuration Free

The BGP-PIC solution depends on internal structures and procedures and does not require any configuration and operator involvement.

7.5. Incremental Deployment

As soon as one router supports BGP-PIC solution, it is possible to benefit from all its benefits (most notably convergence that does not depend in the number of prefixes) without any requirement for other routers to support BGP-PIC.

8. Security Considerations

The behavior described in this document is internal functionality to a router that result in significant improvement to convergence time as well as reduction in CPU and memory used by FIB while not showing change in basic routing and forwarding functionality. As such no additional security risk is introduced by using the mechanisms described in this document.

9. IANA Considerations

This document has no IANA actions.

10. References

10.1. Normative References

- [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>>.
- [RFC3031] Rosen, E., Viswanathan, A., and R. Callon, "Multiprotocol Label Switching Architecture", RFC 3031, DOI 10.17487/RFC3031, January 2001, <<https://www.rfc-editor.org/info/rfc3031>>.

10.2. Informative References

- [I-D.ietf-idr-best-external] Marques, P., Fernando, R., Chen, E., Mohapatra, P., and H. Gredler, "Advertisement of the best external route in BGP", Work in Progress, Internet-Draft, draft-ietf-idr-best-external-05, 3 January 2012, <<https://datatracker.ietf.org/doc/html/draft-ietf-idr-best-external-05>>.
- [RFC5565] Wu, J., Cui, Y., Metz, C., and E. Rosen, "Softwire Mesh Framework", RFC 5565, DOI 10.17487/RFC5565, June 2009, <<https://www.rfc-editor.org/info/rfc5565>>.
- [RFC4364] Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)", RFC 4364, DOI 10.17487/RFC4364, February 2006, <<https://www.rfc-editor.org/info/rfc4364>>.

- [RFC4798] De Clercq, J., Ooms, D., Prevost, S., and F. Le Faucheur, "Connecting IPv6 Islands over IPv4 MPLS Using IPv6 Provider Edge Routers (6PE)", RFC 4798, DOI 10.17487/RFC4798, February 2007, <<https://www.rfc-editor.org/info/rfc4798>>.
- [RFC5036] Andersson, L., Ed., Minei, I., Ed., and B. Thomas, Ed., "LDP Specification", RFC 5036, DOI 10.17487/RFC5036, October 2007, <<https://www.rfc-editor.org/info/rfc5036>>.
- [RFC7911] Walton, D., Retana, A., Chen, E., and J. Scudder, "Advertisement of Multiple Paths in BGP", RFC 7911, DOI 10.17487/RFC7911, July 2016, <<https://www.rfc-editor.org/info/rfc7911>>.
- [RFC6774] Raszuk, R., Ed., Fernando, R., Patel, K., McPherson, D., and K. Kumaki, "Distribution of Diverse BGP Paths", RFC 6774, DOI 10.17487/RFC6774, November 2012, <<https://www.rfc-editor.org/info/rfc6774>>.
- [I-D.pmohapat-idr-fast-conn-restore]
Mohapatra, P., Fernando, R., Filsfils, C., and R. Raszuk, "Fast Connectivity Restoration Using BGP Add-path", Work in Progress, Internet-Draft, draft-pmohapat-idr-fast-conn-restore-03, 22 January 2013, <<https://datatracker.ietf.org/doc/html/draft-pmohapat-idr-fast-conn-restore-03>>.
- [I-D.ietf-rtgwg-segment-routing-ti-lfa]
Bashandy, A., Litkowski, S., Filsfils, C., Francois, P., Decraene, B., and D. Voyer, "Topology Independent Fast Reroute using Segment Routing", Work in Progress, Internet-Draft, draft-ietf-rtgwg-segment-routing-ti-lfa-21, 12 February 2025, <<https://datatracker.ietf.org/doc/html/draft-ietf-rtgwg-segment-routing-ti-lfa-21>>.
- [RFC5714] Shand, M. and S. Bryant, "IP Fast Reroute Framework", RFC 5714, DOI 10.17487/RFC5714, January 2010, <<https://www.rfc-editor.org/info/rfc5714>>.
- [RFC7490] Bryant, S., Filsfils, C., Previdi, S., Shand, M., and N. So, "Remote Loop-Free Alternate (LFA) Fast Reroute (FRR)", RFC 7490, DOI 10.17487/RFC7490, April 2015, <<https://www.rfc-editor.org/info/rfc7490>>.

- [RFC7812] Atlas, A., Bowers, C., and G. Enyedi, "An Architecture for IP/LDP Fast Reroute Using Maximally Redundant Trees (MRT-FRR)", RFC 7812, DOI 10.17487/RFC7812, June 2016, <<https://www.rfc-editor.org/info/rfc7812>>.
- [RFC8277] Rosen, E., "Using BGP to Bind MPLS Labels to Address Prefixes", RFC 8277, DOI 10.17487/RFC8277, October 2017, <<https://www.rfc-editor.org/info/rfc8277>>.
- [RFC8660] Bashandy, A., Ed., Filsfils, C., Ed., Previdi, S., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing with the MPLS Data Plane", RFC 8660, DOI 10.17487/RFC8660, December 2019, <<https://www.rfc-editor.org/info/rfc8660>>.
- [RFC9107] Raszuk, R., Ed., Decraene, B., Ed., Cassar, C., man, E., and K. Wang, "BGP Optimal Route Reflection (BGP ORR)", RFC 9107, DOI 10.17487/RFC9107, August 2021, <<https://www.rfc-editor.org/info/rfc9107>>.

Appendix A. Acknowledgments

Special thanks to Neeraj Malhotra and Yuri Tsier for the valuable help

Special thanks to Bruno Decraene, Theresa Enghardt, Ines Robles, Luc Andre Burdet, and Alvaro Retana for the valuable comments

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

Appendix B. Handling Platforms with Limited Levels of Hierarchy

This section provides additional details on how to handle platforms with limited number of hierarchical levels.

Let's consider a pathlist associated with the leaf "R1" consisting of the list of paths <P1, P2, ..., Pn>. Assume that the leaf "R1" has an OutLabel-list <L1, L2, ..., Ln>. Suppose the path Pi is a recursive path that resolves via a prefix represented by the leaf "R2". The leaf "R2" itself is pointing to a pathlist consisting of the paths <Q1, Q2, ..., Qm>.

If the platform supports the number of hierarchy levels of the forwarding chain, then a packet that uses the path "Pi" will be forwarded according to the steps in Section 4.

Suppose the platform cannot support the number of hierarchy levels in the forwarding chain. FIB manager needs to reduce the number of hierarchy levels when programming the forwarding chain in the FIB. The idea of reducing the number of hierarchy levels is to "flatten" two chain levels into a single level. The "flattening" steps are as follows

1. FIB manager walks to the parent of "Pi", which is the leaf "R2".
 2. FIB manager extracts the parent pathlist of the leaf "R2", which is <Q1, Q2,..., Qm>.
 3. FIB manager also extracts the OutLabel-list of R2 associated with the leaf "R2". Remember that the OutLabel-list of R2 is <L1, L2,..., Lm>.
 4. FIB manager replaces the path "Pi", with the list of pic- paths <Q1, Q2,..., Qm>.
 5. Hence the path list <P1, P2,..., Pn> now becomes "<P1, P2,...,Pi-1, Q1, Q2,..., Qm, Pi+1, Pn>".
1. The path-index stored inside the locations "Q1", "Q2", ..., "Qm" must all be "i" because the index "i" refers to the label "Li" associated with leaf "R1".
 2. FIB manager attaches an OutLabel-list with the new pathlist as follows: <Unlabeled,..., Unlabeled, L1, L2,..., Lm, Unlabeled, ..., Unlabeled>. The size of the label list associated with the flattened pathlist equals the size of the pathlist. Thus there is a 1-1 mapping between every path in the "flattened" pathlist and the OutLabel-list associated with it.

It is noteworthy to mention that the labels in the OutLabel-list associated with the "flattened" pathlist may be stored in the same memory location as the path itself to avoid additional memory access.

The same steps can be applied to all paths in the pathlist <P1, P2,..., Pn> so that all paths are "flattened" thereby reducing the number of hierarchical levels by one. Note that that "flattening" a pathlist pulls in all paths of the parent pic- paths, a desired feature to utilize all paths at all levels. A platform that has a limit on the number of paths in a pathlist for any given leaf may choose to reduce the number paths using methods that are beyond the scope of this document.

The steps can be recursively applied to other paths at the same levels or other levels to recursively reduce the number of hierarchical levels to an arbitrary value so as to accommodate the capability of the forwarding engine.

Because a flattened pathlist may have an associated OutLabel-list the forwarding behavior has to be slightly modified. The modification is done by adding the following step right after step 4 in Section 4.

1. If there is an OutLabel-list associated with the pathlist, then if the path "Pi" is chosen by the hashing algorithm, retrieve the label at location "i" in that OutLabel-list and apply the label action of that label on the packet.

The steps in this Section to are applied to an example in the next Section.

Appendix C. Example: Flattening a forwarding chain.

This example uses a case of inter-AS option C [RFC4364] where there are 3 levels of hierarchy. Figure 4 illustrates the sample topology. The Autonomous System Border Routers (ASBRs) on the ingress domain (Domain 1) use BGP to advertise the core routers (ASBRs and ePEs) of the egress domain (Domain 2) to the iPE. The end result is that the ingress PE (iPE) has 2 levels of recursion for the VPN prefixes VPN-IP1 and VPN-IP2.

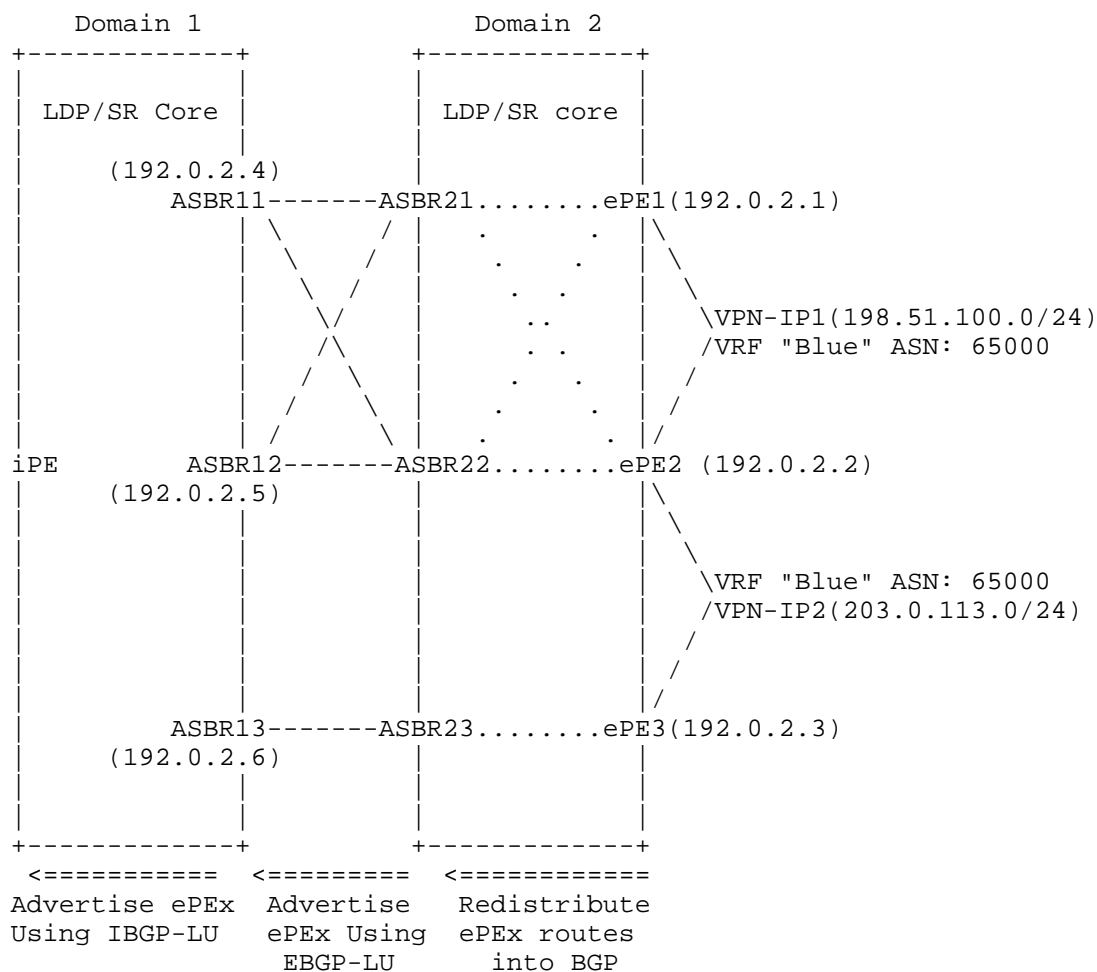


Figure 3: Sample 3-level hierarchy topology

The following assumptions about connectivity are made:

- * In "Domain 2", both ASBR21 and ASBR22 can reach both ePE1 and ePE2 using the same metric.
- * In "Domain 2", only ASBR23 can reach ePE3.
- * In "Domain 1", iPE (the ingress PE) can reach ASBR11, ASBR12, and ASBR13 via IGP using the same metric.

The following assumptions are made about the labels:

- * The VPN labels advertised by ePE1 and ePE2 for prefix VPN-IP1 are VPN-L11 and VPN-L21, respectively.
- * The VPN labels advertised by ePE2 and ePE3 for prefix VPN-IP2 are VPN-L22 and VPN-L32, respectively.
- * The labels advertised by ASBR11 to iPE using BGP-LU for the egress PEs ePE1 and ePE2 are LASBR111(ePE1) and LASBR112(ePE2), respectively.
- * The labels advertised by ASBR12 to iPE using BGP-LU for the egress PEs ePE1 and ePE2 are LASBR121(ePE1) and LASBR122(ePE2), respectively.
- * The label advertised by ASBR13 to iPE using BGP-LU for the egress PE ePE3 is LASBR13(ePE3).
- * The IGP labels advertised by the next hops directly connected to iPE towards ASBR11, ASBR12, and ASBR13 in the core of domain 1 are IGP-L11, IGP-L12, and IGP-L13, respectively.
- * Both the routers ASBR21 and ASBR22 of Domain 2 advertise the same label LASBR21 and LASBR22 for the egress PEs ePE1 and ePE2, respectively, to the routers ASBR11 and ASBR22 of Domain 1.
- * The router ASBR23 of Domain 2 advertises the label LASBR23 for the egress PE ePE3 to the router ASBR13 of Domain 1.

Based on these connectivity assumptions and the topology in Figure 4, the routing table on iPE is

```
65000: 198.51.100.0/24
    via ePE1 (192.0.2.1), VPN Label: VPN-L11
    via ePE2 (192.0.2.2), VPN Label: VPN-L21
65000: 203.0.113.0/24
    via ePE2 (192.0.2.2), VPN Label: VPN-L22
    via ePE3 (192.0.2.3), VPN Label: VPN-L32

192.0.2.1/32 (ePE1)
    via ASBR11, Label: LASBR111(ePE1) via ASBR12, Label:
    LASBR121(ePE1)

192.0.2.2/32 (ePE2)
    via ASBR11, Label: LASBR112(ePE2) via ASBR12, Label:
    LASBR122(ePE2)

192.0.2.3/32 (ePE3)
    Via ASBR13, Label: LASBR13(ePE3)
```

```
192.0.2.4/32 (ASBR11)
  via Core, Label:    IGP-L11
192.0.2.5/32 (ASBR12)
  via Core, Label:    IGP-L12
192.0.2.6/32 (ASBR13)
  via Core, Label:    IGP-L13
```

The diagram in Figure 5 illustrates the forwarding chain in iPE assuming that the forwarding hardware in iPE supports 3 levels of hierarchy. The leaves corresponding to the ASBRs on domain 1 (ASBR11, ASBR12, and ASBR13) are at the bottom of the hierarchy. There are few important points:

- * Because the hardware supports the required depth of hierarchy, the sizes of a pathlist equal the size of the label list associated with the leaves using this pathlist.
- * The path-index inside the pathlist entry indicates the label that will be picked from the OutLabel-List associated with the child leaf if that path is chosen by the forwarding engine hashing function.

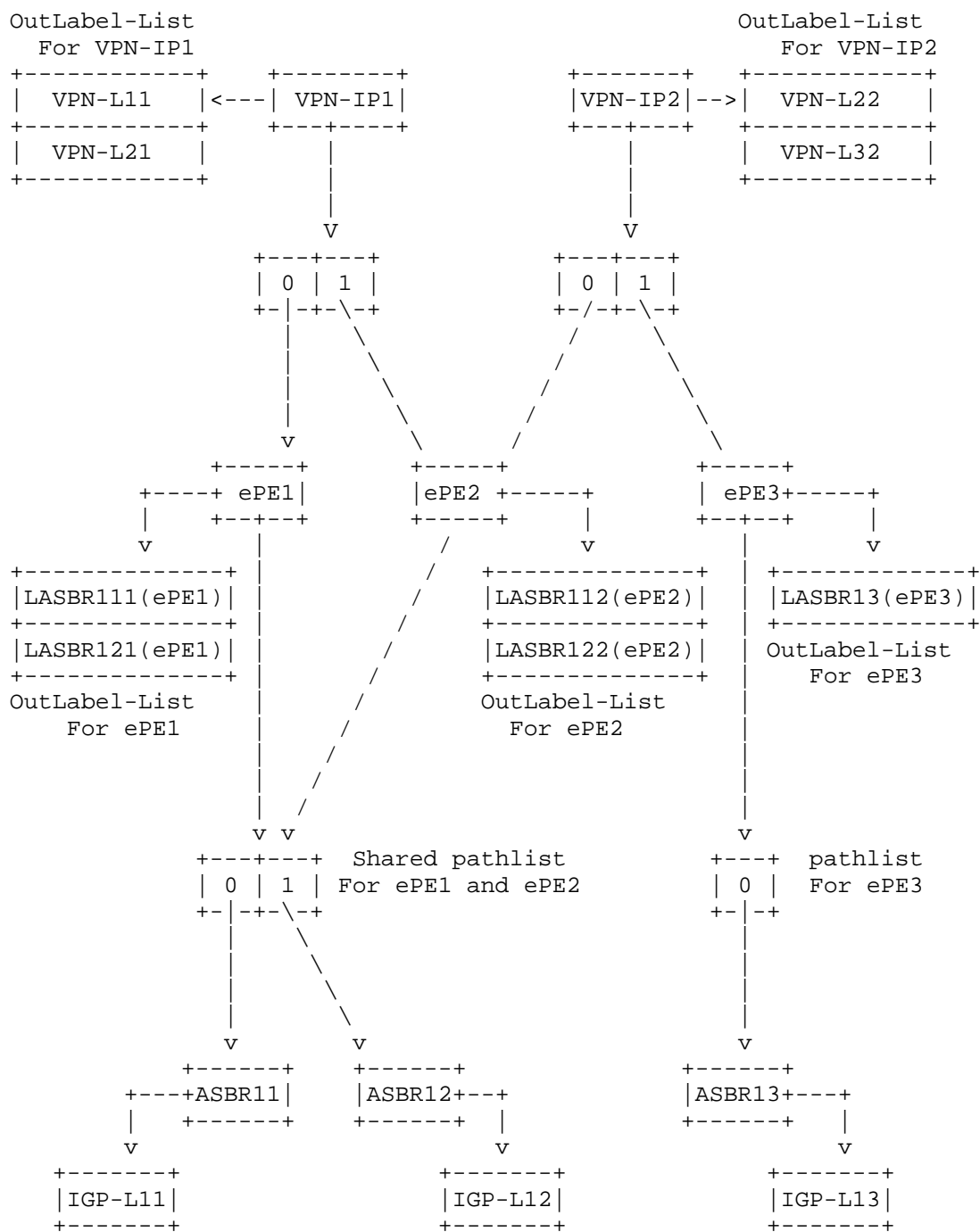


Figure 4: Forwarding Chain for hardware supporting 3 Levels

Now suppose the hardware on iPE (the ingress PE) supports 2 levels of hierarchy only. In that case, the 3-levels forwarding chain in Figure 5 needs to be "flattened" into 2 levels only.

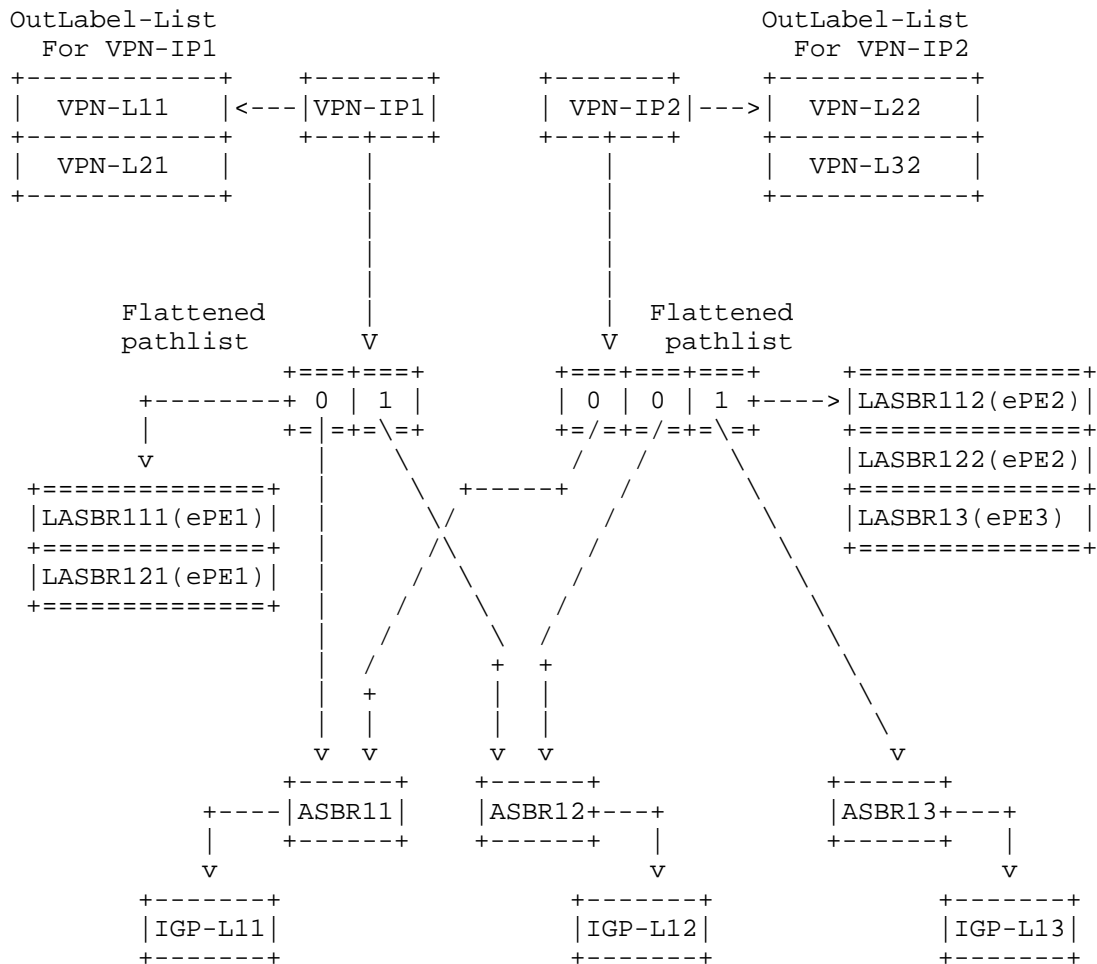


Figure 5: Flattening 3 levels to 2 levels of Hierarchy on iPE

Figure 6 represents one way to "flatten" a 3 levels hierarchy into two levels. There are a few important points:

- * As mentioned in Appendix B, a flattened pathlist may have label lists associated with them. The size of the label list associated with a flattened pathlist equals the size of the pathlist. Hence it is possible that an implementation includes these label lists in the flattened pathlist itself.
- * Again as mentioned in Appendix B, the size of a flattened pathlist may not be equal to the size of the OutLabel-lists of leaves using the flattened pathlist. So the indices inside a flattened pathlist still indicate the label index in the OutLabel-Lists of the leaves using that pathlist. Because the size of the flattened pathlist may be different from the size of the OutLabel-lists of the leaves, the indices may be repeated.
- * Let's take a look at the flattened pathlist used by the prefix "VPN-IP2". The pathlist associated with the prefix "VPN-IP2" has three entries.
 - The first and second entry have index "0". This is because both entries correspond to ePE2. Thus when hashing performed by the forwarding engine results in using the first or the second entry in the pathlist, the forwarding engine will pick the correct VPN label "VPN-L22", which is the label advertised by ePE2 for the prefix "VPN-IP2".
 - The third entry has the index "1". This is because the third entry corresponds to ePE3. Thus when the hashing is performed by the forwarding engine results in using the third entry in the flattened pathlist, the forwarding engine will pick the correct VPN label "VPN-L32", which is the label advertised by "ePE3" for the prefix "VPN-IP2".

Now let's try and apply the forwarding steps in Section 4 together with the additional step in Section Appendix B to the flattened forwarding chain illustrated in Figure 6.

- * Suppose a packet arrives at "iPE" and matches the VPN prefix "VPN-IP2".
- * The forwarding engine walks to the parent of the "VPN-IP2", which is the flattened pathlist and applies a hashing algorithm to pick a path.
- * Suppose the hashing by the forwarding engine picks the second path in the flattened pathlist associated with the leaf "VPN-IP2".
- * Because the second path has the index "0", the label "VPN-L22" is pushed on the packet.

- * Next the forwarding engine picks the second label from the OutLabel-List associated with the flattened pathlist resulting in "LASBR122(ePE2)" being the next pushed label.
- * The forwarding engine now moves to the parent of the flattened pathlist corresponding to the second path. The parent is the IGP label leaf corresponding to "ASBR12".
- * So the packet is forwarded towards the ASBR "ASBR12" and the IGP label at the top will be "IGP-L12".

Based on the above steps, a packet arriving at iPE and destined to the prefix VPN-L22 reaches its destination as follows:

- o iPE sends the packet along the shortest path towards ASBR12 with the following label stack starting from the top: {L12, LASBR122(ePE2), VPN-L22}.
- o The penultimate hop of ASBR12 pops the top label "L12". Hence the packet arrives at ASBR12 with the remaining label stack {LASBR122(ePE2), VPN-L22} where "LASBR122(ePE2)" is the top label.
- o ASBR12 swaps "LASBR122(ePE2)" with the label "LASBR22(ePE2)", which is the label advertised by ASBR22 for the ePE2 (the egress PE).
- o ASBR22 receives the packet with "LASBR22(ePE2)" at the top.
- o Hence ASBR22 swaps "LASBR22(ePE2)" with the IGP label for ePE2 advertised by the next-hop towards ePE2 in domain 2, and sends the packet along the shortest path towards ePE2.
- o The penultimate hop of ePE2 pops the top label. Hence ePE2 receives the packet with the top label VPN-L22 at the top
- o ePE2 pops "VPN-L22" and sends the packet as a pure IP packet towards the destination VPN-IP2.

Appendix D. Perspective

The following table puts the BGP-PIC benefits in perspective assuming

- * 1M impacted BGP prefixes
- * IGP convergence ~ 500 msec
- * local protection ~ 50msec

- * FIB Update per BGP destination ~ 100usec conservative,
~ 10usec optimistic

- * BGP best route recalculation per BGP destination
~ 10usec optimistic,
~ 100usec optimistic

Without PIC With PIC

Local IGP Failure 10 to 100sec 50msec

Local BGP Failure 100 to 200sec 50msec

Remote IGP Failure 10 to 100sec 500msec

Local BGP Failure 100 to 200sec 500msec

Upon local IGP next-hop failure or remote IGP next-hop failure, the existing primary BGP next-hop is intact and usable hence the resiliency only depends on the ability of the FIB mechanism to reflect the new path to the BGP next-hop to the depending BGP destinations. Without BGP-PIC, a conservative back-of-the-envelope estimation for this FIB update is 100usec per BGP destination. An optimistic estimation is 10usec per entry.

Upon local BGP next-hop failure or remote BGP next-hop failure, without the BGP-PIC mechanism, a new BGP Best-Path needs to be recomputed and new updates need to be sent to peers. This depends on BGP processing time that will be shared between best-path computation, RIB update and peer update. A conservative back-of-the-envelope estimation for this is 200usec per BGP destination. An optimistic estimation is 100usec per entry.

Authors' Addresses

Ahmed Bashandy (editor)
HPE
United States of America
Email: abashandy.ietf@gmail.com

Clarence Filsfils
Cisco Systems
Email: cfilsfil@cisco.com

Pradosh Mohapatra
Sproute Networks
United States of America
Email: mpradosh@yahoo.com

Yingzhen Qu (editor)
Futurewei Technologies
United States of America
Email: yingzhen.ietf@gmail.com