

LSR Working Group  
Internet Draft  
Intended status: Standards Track  
Expires: April 2019

Dave Allan  
Ericsson  
October 2018

A Distributed Algorithm for Constrained Flooding of IGP  
Advertisements  
draft-allan-lsr-flooding-algorithm-00

Abstract

This document describes a distributed algorithm that can be applied to the problem of constraining IGP flooding in dense mesh topologies. The flooding topology utilizes two node-diverse spanning trees in order to provide complete coverage in the presence of any single failure while constraining the number of LSAs received by any IGP speaker connected to the flooding topology.

Status of this Memo

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

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/lid-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire in March 2019.

Copyright and License Notice

Copyright (c) 2018 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

1. Introduction.....	3
1.1. Authors.....	3
1.2. Requirements Language.....	3
2. Conventions used in this document.....	3
2.1. Terminology.....	3
3. Solution Overview.....	4
3.1. The Flooding Topology.....	4
3.2. Solution Applicability.....	4
3.3. Algorithm.....	4
3.3.1. Algorithm Basics.....	5
3.3.2. Generating Diverse Trees.....	5
3.3.3. Desirable Properties Computation Wise.....	6
4. Applying the Algorithm.....	6
4.1. Tree Generation.....	6
4.2. Illustrating the result.....	6
4.3. Interactions between Participating and Non-Participating Nodes.....	7
4.4. Flooding of LSAs.....	8
4.5. Root Selection.....	9
4.6. Node Additions.....	9
5. Further work.....	10
5.1. Thoughts on Coexistence in the Context of a Larger Network..	10
5.1.1. Multiple flooding Domains and the Severing of Flooding Domains.....	10
5.2. Thoughts on Flooding Topology Re-Optimization.....	10
5.3. Thoughts on Node and Network Initialization.....	11
5.4. Thoughts on Loop Prevention.....	11
5.5. Thoughts on Pathological Failure Scenarios.....	11
6. Acknowledgements.....	12
7. Security Considerations.....	12
8. IANA Considerations.....	12
9. References.....	12

9.1. Normative References.....	12
9.2. Informative References.....	12
10. Author's Address.....	13

## 1. Introduction

This memo describes an algorithm suitable for reducing the quantity of IGP flooding in dense mesh networks. The only property that the algorithm is dependent upon is that there are at least two equal and diverse shortest paths between any pair of IGP speakers in order to meet the requirements elucidated in [Li]. The algorithm uses a re-purposing of the tie breaking algorithm used in 802.1aq Shortest Path Bridging as an element of construction of the flooding topology. It is not the intention of this memo to specify a complete solution, but to offer a foundation of an eventual solution.

### 1.1. Authors

David Allan

### 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 RFC2119 [RFC2119].

## 2. Conventions used in this document

### 2.1. Terminology

Member Adjacency - An adjacency that has been determined to part of the flooding topology.

Member Node - A Participant node that is connected to the flooding topology.

Participant Adjacency - An adjacency between two participating nodes. It may be a member adjacency or a non-member adjacency

Non-Participant Adjacency - An adjacency where at least one of the two nodes is a not a Participating Node

Participating Node - An IGP speaker that has advertised the capability, and hence the intention, to participate in a flooding topology

### 3. Solution Overview

#### 3.1. The Flooding Topology

A flooding topology is composed of a contiguously connected set of participating nodes.

The flooding topology constructed from two diversely rooted spanning trees. A participating node that is connected to the physical topology with a degree of two or greater and has at least two participating adjacencies will be bi-connected to the flooding topology.

The resulting flooding topology diameter will typically be two times the depth of the tree hierarchy. The compromise in this approach is that a subset of nodes in the network will not see a reduction of the replication burden from current practice when flooding LSAs as the degree of a subset of nodes in the flooding topology will correspond to the degree of the physical topology.

The protocol structure of flooded information is unmodified. A participant node may relay a received LSA onto member links of both spanning trees. Specific forwarding rules prevent undue flooding, the result being that every participant node that is bi-connected to the flooding topology will receive two copies of any flooded LSA in a fault free network. Participating nodes that due to network degradation are only singly connected will receive one copy. The forwarding rules are described in section 4.4.

#### 3.2. Solution Applicability

This algorithm has been considered in the context of pure bipartite graphs, bipartite graphs modified with the addition of intra-tier adjacencies, and hierarchical variations of the above. Applicability to other network designs is for further study.

For all graphs the link costs are assumed to be common for all inter-tier links and common for any intra-tier links. Inter-tier and intra-tier links do not have to have the same cost.

#### 3.3. Algorithm

The algorithm borrows from 802.1aq for the construction of the spanning trees used in this application. This is described in clause 28.5 of [802.1Q].

### 3.3.1. Algorithm Basics

The key component of the 802.1aq employed is the tie breaking algorithm. The original application of the algorithm was to produce a symmetrically congruent mesh of multicast trees and unicast forwarding whereby the path between any two nodes in the network was symmetric in both directions and congruent for both unicast and multicast traffic.

For this application the algorithm is used in the generation of two diversely rooted spanning trees that define the flooding topology.

As part of tree construction, the algorithm tie breaks between equal cost paths. When a tie is identified as part of a Dijkstra computation, a path-id is constructed for each equal cost path. A path-id is expressed as a lexicographically sorted list of the node-ids in the path. The set of equal cost paths is ranked, and the lowest selected. As an example:

Path-id 23-39-44-68-85 is ranked lower than

Path-id 23-44-59-63-90

When the path-ids are of unequal length, the path-ids with the fewest hops are ranked superior to the longer paths, and tie breaking is applied to select between the shorter path-ids. This is not expected to apply in the general case of the dense graphs this application is targeted at.

The node-ids used would be the loopback address of each node, therefore each path-id will be unique.

### 3.3.2. Generating Diverse Trees

The algorithm includes the concept of an "algorithm-mask", which is a value XOR'd with the node-ids prior to sorting into path IDs and ranking the paths. This permits the construction of diverse trees in a dense topology.

Two algorithm masks are used (zero and -1). When computing two trees from the same root, when there are at least two nodes to choose from at each distance from the root, fully diverse trees will be generated. When computing two trees from diverse roots in a tree architecture, diverse nodes will be selected in each tier in the hierarchy as the relay nodes to the next tier.

### 3.3.3. Desirable Properties Computation Wise

The algorithm has the property of permitting the pruning of intermediate state as a Dijkstra progresses as ties can be immediately evaluated, and the all but the selected path removed from further consideration. This is desirable when computing a Dijkstra in a dense graph as all path permutations do not need to be carried forward during computation. This permits the computation to be quite fast.

The resulting computational complexity would still be expressed as  $2N(\ln N)$ .

## 4. Applying the Algorithm

### 4.1. Tree Generation

Each IGP speaker in the network has knowledge of each of the two spanning tree roots and the algorithm mask associated with each. This memo does not specify how root selection is performed and disseminated through the network, but does discuss selection requirements in section 4.5.

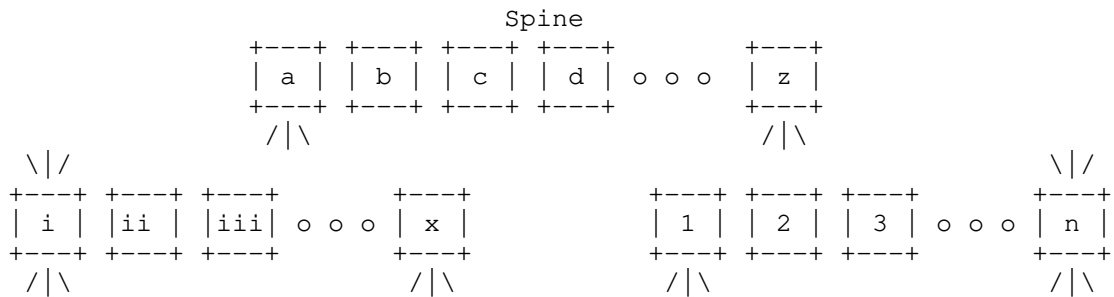
Each root has one of the two algorithm masks associated with it.

Each participating IGP speaker in the network computes a spanning tree from each the two roots (using the algorithm mask associated with each root) and from that can determine its own role in the flooding topology. The two spanning trees are designated the "low spanning tree" and the "high spanning tree".

The spanning trees are a starting point for a redundant topology. Unlike the commonly accepted operation of a spanning tree, in this application the distinction between upstream and downstream adjacencies is important and is an input to how a member node further relays any LSAs received. Upstream member adjacencies are in the direction of a root, and downstream member adjacencies are in the direction away from the root.

### 4.2. Illustrating the result

The following diagram illustrates the general layout of the flooding graph constructed using the algorithm as applied to a bi-partite style of tree (no intra tier links):



In the example, there are two tiers of switches. The spine (nodes a..z), and the next tier with two groups of nodes (i..x) and (1..n). The algorithm will select the node with the lowest node ID in each tier as the replicating node for the low spanning tree; 'a' and 'i' for the set of nodes connecting the spine and the next tier. The algorithm will select the nodes with highest node ID in the same set of nodes for the high spanning tree; 'z' and 'n' for the same set of nodes.

In the flooding topology:

- Node 'a' is connected to nodes i..x and 1..n for the low spanning tree.
- Node 'z' is connected to the same set of nodes for the high spanning tree.
- Node 'i' is connected to nodes 'a'..'z' for the low spanning tree, and
- Node 'n' is connected to the same nodes for the high spanning tree.
- All other nodes are bi-connected to the flooding topology

If there was a further tier added below nodes i..x, then 'i' and 'x' would be selected as the replicating nodes for the low and high spanning tree respectively. This is similarly true for nodes 1..n.

#### 4.3. Interactions between Participating and Non-Participating Nodes

This solution proposes primarily only nodal behaviors with respect to constraining flooding to member adjacencies. To address the scenario

where the participating nodes were a subset of a larger network, it would be necessary to advertise the capability to participate in flood reduction.

This would then require that each participating node use this information to be able to identify the set of participating adjacencies and confine the spanning tree computation to the set of participating adjacencies in order to identify local set of member adjacencies. Interactions with non-participant adjacencies would conform to current practice.

#### 4.4. Flooding of LSAs

The design of the protocol elements that are flooded is unmodified by this solution. Therefore, there is no additional information available to associate a received LSA with a given tree, nor is such information needed; the two spanning trees are not treated as unique entities in the flooding topology.

As per current practice, a node does not relay LSAs that it has already seen.

A new LSA received from an upstream member adjacency is flooded on:

- All downstream member adjacencies exclusive of the adjacency of arrival, irrespective of which tree the adjacencies are part of.
- All non-participant adjacencies

A new LSA received from a downstream member adjacency is flooded on:

- All other member adjacencies exclusive of the adjacency of arrival irrespective of which tree the adjacencies are part of.
- All non-participant adjacencies

A new LSA received from a member adjacency where upstream and downstream is ambiguous (it is an upstream member on one of the spanning trees and a downstream member on the other), is flooded on:

- All other member adjacencies exclusive of the adjacency of arrival irrespective of which adjacency the links are part of.
- All Non-Participant adjacencies



A new LSA received from a non-member adjacency is flooded on all member adjacency irrespective of which tree the adjacencies are part of (see sections 5.1 and 5.5).

#### 4.5. Root Selection

The algorithm depends on tie breaking between sets of node IDs to produce diverse paths, therefore it does place some restrictions on root selection.

A root SHOULD be selected so that the root's node-id when XORd with the associated algorithm mask is the lowest ranked node in the local tier in the tree hierarchy. This would be analogous to path-id ranking where the paths were all of length 1.

The root MUST NOT be selected such that the node-ID when XORd with the other root's algorithm mask is the lowest ranked node. This would result in the root also being a transit node for the other spanning tree and produce a scenario whereby a single failure could render both spanning trees incomplete.

Roots MUST NOT be directly connected for either of the low or high spanning trees. If the topology does not permit this to be satisfied purely by root selection, then the inter-root adjacency must be pruned from the graph prior to spanning tree computation to ensure that diverse paths between the roots are used.

For a true bipartite graph, there are no other restrictions on node selection.

For a bipartite graph modified with inter-tier links, the roots MUST be placed in different tiers to ensure a pathological combination of link weights and node-ids does not result in a scenario where a single failure would render the flooding topology incomplete.

Other sources of failure may exist that may require an administrative component to root selection. This, for example, would ensure that both roots were not selected from a common shared risk group.

See also section 5.5.

#### 4.6. Node Additions

A participating node that is added to the topology will initially not be served by the flooding topology. A participating node adjacent to that node is required to treat it as a non-participating node until such time as tree re-optimization has completed. At the end of tree

optimization, typically two adjacent participating nodes will have member adjacencies with the new node, so the ability to flood LSAs between the new node and the flooding topology will have been uninterrupted during the process.

## 5. Further work

### 5.1. Thoughts on Coexistence in the Context of a Larger Network

A node that had a combination of participating and non-participating adjacencies would be required to do the following:

- For any new LSA received on a participating adjacency, in addition to the rules for member adjacencies, it would also flood the LSA on all non-participating adjacencies.
- For any new LSA received on a non-participating adjacency, it would flood the LSA on all member adjacencies.

This is reflected in the forwarding rules described in section 4.4.

#### 5.1.1. Multiple flooding Domains and the Severing of Flooding Domains

It is possible to envision several scenarios whereby there are sets of participating nodes that are not contiguously connected via participating adjacencies in a given IGP domain.

1. A node has been incorrectly configured as a participating node but has no participating adjacencies.
2. A participating node or set of nodes has become severed from the flooding topology but is still connected to other nodes in the network. Nodes in this set would still be able to compute a local extension of the flooding topology, but it would only be useful if the set was sufficiently large that a majority of the nodes were not connected to non-participants.
3. Procedures are designed to permit more than one flooding topology in an IGP domain. In which case participating nodes would have to be administratively configured to associate with a flooding topology instance.

### 5.2. Thoughts on Flooding Topology Re-Optimization

After a topology change, it is desirable that the flooding topology remain stable until the network has stabilized. However a single failure may render one of the spanning trees incomplete, such that a

further single failure could make the flooding topology incomplete. Therefore procedures should include re-optimization of the flooding topology after a topology change. In order to maintain complete coverage it would make sense not to recompute the spanning trees simultaneously.

One approach that would appear to make sense to separate in time network convergence, re-optimization of the low spanning tree and re-optimization of the high spanning tree.

The ideal would be to reoptimize an incomplete tree first, however this would require the participating nodes to maintain a complete map of all member adjacencies so that a common determination of the most degraded spanning tree and hence the order of re-optimization could be made.

### 5.3. Thoughts on Node and Network Initialization

A participating node at power up will be not be able to establish member links until it has synchronized with the network and the network is stable in the new topology. This suggests it simply treats power up similarly to how a topology change and network re-optimization is treated. The only difference being that it will flood all LSAs received or originated as per current practice until both spanning trees have stabilized.

### 5.4. Thoughts on Loop Prevention

802.1aq included additional mechanisms to prevent looping, a reverse path forwarding check, and digest exchange across adjacencies to ensure IGP synchronization.

Routing LSAs are not relayed if they are a duplicate, therefore destructive looping cannot occur and additional mitigation mechanisms are not required.

### 5.5. Thoughts on Pathological Failure Scenarios

While in a stable fault free network with sufficient mesh density of the types considered, the flooding topology used by this solution would ensure that no single failure rendered both spanning trees incomplete, it is also useful to consider multiple failure scenarios and if they can be mitigated.

Preliminary analysis suggests that in a tree network of sufficient mesh density, the only dual link failure that can render the flooding topology incomplete is if a participant node has failures in both

upstream member adjacencies. This can be partially mitigated if the node recognizes this scenario and reverts to flooding on all adjacencies. If the suggested procedures of 5.1.1 above are adopted, surrounding participating nodes that receive the LSA on a non-member adjacency will introduce the LSA into the flooding topology.

The pathological scenario is the simultaneous failure of both roots. This does suggest that root selection should place the roots two hops apart so there will be a constituency of participants that would observe a simultaneous failure of both upstream member adjacencies and revert to normal flooding.

## 6. Acknowledgements

The author would like to acknowledge Jerome Chiabaut for his original algorithm work that underpins this memo.

## 7. Security Considerations

For a future version of this document.

## 8. IANA Considerations

This memo requires no IANA allocations

## 9. References

### 9.1. Normative References

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.

### 9.2. Informative References

[802.1Q] 802.1Q (2014) IEEE Standard for Local and Metropolitan Area Networks--Media Access Control (MAC) Bridges and Virtual Bridged Local Area Networks

[Li] Li, T., Psenak, P., "Dynamic Flooding on Dense Graphs", IETF work in progress, draft-li-dynamic-flooding-05, June 2018

10. Author's Address

Dave Allan  
Ericsson  
2455 Augustine Drive  
Santa Clara, CA 95054  
USA  
Email: david.i.allan@ericsson.com

RIFT Working Group  
Internet-Draft  
Intended status: Standards Track  
Expires: November 27, 2020

A. Przygienda, Ed.  
Juniper  
A. Sharma  
Comcast  
P. Thubert  
Cisco  
Bruno. Rijsman  
Individual  
Dmitry. Afanasiev  
Yandex  
May 26, 2020

RIFT: Routing in Fat Trees  
draft-ietf-rift-rift-12

Abstract

This document defines a specialized, dynamic routing protocol for Clos and fat-tree network topologies optimized towards minimization of configuration and operational complexity. The protocol

- o deals with no configuration, fully automated construction of fat-tree topologies based on detection of links,
- o minimizes the amount of routing state held at each level,
- o automatically prunes and load balances topology flooding exchanges over a sufficient subset of links,
- o supports automatic disaggregation of prefixes on link and node failures to prevent black-holing and suboptimal routing,
- o allows traffic steering and re-routing policies,
- o allows loop-free non-ECMP forwarding,
- o automatically re-balances traffic towards the spines based on bandwidth available and finally
- o provides mechanisms to synchronize a limited key-value data-store that can be used after protocol convergence to e.g. bootstrap higher levels of functionality on nodes.

## 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 November 27, 2020.

## Copyright Notice

Copyright (c) 2020 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 Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

1. Authors . . . . .	6
2. Introduction . . . . .	6
2.1. Requirements Language . . . . .	8
3. Reference Frame . . . . .	8
3.1. Terminology . . . . .	8
3.2. Topology . . . . .	13
4. RIFT: Routing in Fat Trees . . . . .	15
4.1. Overview . . . . .	16
4.1.1. Properties . . . . .	16
4.1.2. Generalized Topology View . . . . .	17
4.1.2.1. Terminology . . . . .	17
4.1.2.2. Clos as Crossed Crossbars . . . . .	18
4.1.3. Fallen Leaf Problem . . . . .	28
4.1.4. Discovering Fallen Leaves . . . . .	30

4.1.5.	Addressing the Fallen Leaves Problem . . . . .	31
4.2.	Specification . . . . .	32
4.2.1.	Transport . . . . .	33
4.2.2.	Link (Neighbor) Discovery (LIE Exchange) . . . . .	33
4.2.2.1.	LIE FSM . . . . .	36
4.2.3.	Topology Exchange (TIE Exchange) . . . . .	46
4.2.3.1.	Topology Information Elements . . . . .	46
4.2.3.2.	South- and Northbound Representation . . . . .	46
4.2.3.3.	Flooding . . . . .	49
4.2.3.4.	TIE Flooding Scopes . . . . .	56
4.2.3.5.	'Flood Only Node TIEs' Bit . . . . .	59
4.2.3.6.	Initial and Periodic Database Synchronization . . . . .	60
4.2.3.7.	Purging and Roll-Overs . . . . .	60
4.2.3.8.	Southbound Default Route Origination . . . . .	61
4.2.3.9.	Northbound TIE Flooding Reduction . . . . .	62
4.2.3.10.	Special Considerations . . . . .	67
4.2.4.	Reachability Computation . . . . .	67
4.2.4.1.	Northbound SPF . . . . .	67
4.2.4.2.	Southbound SPF . . . . .	68
4.2.4.3.	East-West Forwarding Within a non-ToF Level . . . . .	69
4.2.4.4.	East-West Links Within ToF Level . . . . .	69
4.2.5.	Automatic Disaggregation on Link & Node Failures . . . . .	69
4.2.5.1.	Positive, Non-transitive Disaggregation . . . . .	69
4.2.5.2.	Negative, Transitive Disaggregation for Fallen Leaves . . . . .	73
4.2.6.	Attaching Prefixes . . . . .	75
4.2.7.	Optional Zero Touch Provisioning (ZTP) . . . . .	84
4.2.7.1.	Terminology . . . . .	85
4.2.7.2.	Automatic SystemID Selection . . . . .	86
4.2.7.3.	Generic Fabric Example . . . . .	87
4.2.7.4.	Level Determination Procedure . . . . .	88
4.2.7.5.	ZTP FSM . . . . .	89
4.2.7.6.	Resulting Topologies . . . . .	95
4.2.8.	Stability Considerations . . . . .	97
4.3.	Further Mechanisms . . . . .	98
4.3.1.	Overload Bit . . . . .	98
4.3.2.	Optimized Route Computation on Leaves . . . . .	98
4.3.3.	Mobility . . . . .	98
4.3.3.1.	Clock Comparison . . . . .	100
4.3.3.2.	Interaction between Time Stamps and Sequence Counters . . . . .	100
4.3.3.3.	Anycast vs. Unicast . . . . .	101
4.3.3.4.	Overlays and Signaling . . . . .	101
4.3.4.	Key/Value Store . . . . .	101
4.3.4.1.	Southbound . . . . .	101
4.3.4.2.	Northbound . . . . .	102
4.3.5.	Interactions with BFD . . . . .	102
4.3.6.	Fabric Bandwidth Balancing . . . . .	103



4.3.6.1.	Northbound Direction . . . . .	103
4.3.6.2.	Southbound Direction . . . . .	106
4.3.7.	Label Binding . . . . .	106
4.3.8.	Leaf to Leaf Procedures . . . . .	106
4.3.9.	Address Family and Multi Topology Considerations . . . . .	106
4.3.10.	Reachability of Internal Nodes in the Fabric . . . . .	107
4.3.11.	One-Hop Healing of Levels with East-West Links . . . . .	107
4.4.	Security . . . . .	107
4.4.1.	Security Model . . . . .	107
4.4.2.	Security Mechanisms . . . . .	109
4.4.3.	Security Envelope . . . . .	110
4.4.4.	Weak Nonces . . . . .	113
4.4.5.	Lifetime . . . . .	114
4.4.6.	Key Management . . . . .	114
4.4.7.	Security Association Changes . . . . .	114
5.	Examples . . . . .	115
5.1.	Normal Operation . . . . .	115
5.2.	Leaf Link Failure . . . . .	117
5.3.	Partitioned Fabric . . . . .	118
5.4.	Northbound Partitioned Router and Optional East-West Links . . . . .	119
6.	Implementation and Operation: Further Details . . . . .	120
6.1.	Considerations for Leaf-Only Implementation . . . . .	120
6.2.	Considerations for Spine Implementation . . . . .	121
6.3.	Adaptations to Other Proposed Data Center Topologies . . . . .	121
6.4.	Originating Non-Default Route Southbound . . . . .	122
7.	Security Considerations . . . . .	122
7.1.	General . . . . .	122
7.2.	ZTP . . . . .	122
7.3.	Lifetime . . . . .	123
7.4.	Packet Number . . . . .	123
7.5.	Outer Fingerprint Attacks . . . . .	123
7.6.	TIE Origin Fingerprint DoS Attacks . . . . .	123
7.7.	Host Implementations . . . . .	124
8.	IANA Considerations . . . . .	124
8.1.	Requested Multicast and Port Numbers . . . . .	124
8.2.	Requested Registries with Suggested Values . . . . .	125
8.2.1.	Registry RIFT_v4/common/AddressFamilyType . . . . .	125
8.2.1.1.	Requested Entries . . . . .	125
8.2.2.	Registry RIFT_v4/common/HierarchyIndications . . . . .	125
8.2.2.1.	Requested Entries . . . . .	125
8.2.3.	Registry RIFT_v4/common/IEEE802_1ASTimeStampType . . . . .	125
8.2.3.1.	Requested Entries . . . . .	126
8.2.4.	Registry RIFT_v4/common/IPAddressType . . . . .	126
8.2.4.1.	Requested Entries . . . . .	126
8.2.5.	Registry RIFT_v4/common/IPPrefixType . . . . .	126
8.2.5.1.	Requested Entries . . . . .	126
8.2.6.	Registry RIFT_v4/common/IPv4PrefixType . . . . .	126

8.2.6.1. Requested Entries . . . . .	126
8.2.7. Registry RIFT_v4/common/IPv6PrefixType . . . . .	126
8.2.7.1. Requested Entries . . . . .	127
8.2.8. Registry RIFT_v4/common/PrefixSequenceType . . . . .	127
8.2.8.1. Requested Entries . . . . .	127
8.2.9. Registry RIFT_v4/common/RouteType . . . . .	127
8.2.9.1. Requested Entries . . . . .	127
8.2.10. Registry RIFT_v4/common/TIETypeType . . . . .	128
8.2.10.1. Requested Entries . . . . .	128
8.2.11. Registry RIFT_v4/common/TieDirectionType . . . . .	128
8.2.11.1. Requested Entries . . . . .	128
8.2.12. Registry RIFT_v4/encoding/Community . . . . .	128
8.2.12.1. Requested Entries . . . . .	128
8.2.13. Registry RIFT_v4/encoding/KeyValueTIEElement . . . . .	129
8.2.13.1. Requested Entries . . . . .	129
8.2.14. Registry RIFT_v4/encoding/LIEPacket . . . . .	129
8.2.14.1. Requested Entries . . . . .	129
8.2.15. Registry RIFT_v4/encoding/LinkCapabilities . . . . .	130
8.2.15.1. Requested Entries . . . . .	130
8.2.16. Registry RIFT_v4/encoding/LinkIDPair . . . . .	130
8.2.16.1. Requested Entries . . . . .	131
8.2.17. Registry RIFT_v4/encoding/Neighbor . . . . .	131
8.2.17.1. Requested Entries . . . . .	131
8.2.18. Registry RIFT_v4/encoding/NodeCapabilities . . . . .	131
8.2.18.1. Requested Entries . . . . .	132
8.2.19. Registry RIFT_v4/encoding/NodeFlags . . . . .	132
8.2.19.1. Requested Entries . . . . .	132
8.2.20. Registry RIFT_v4/encoding/NodeNeighborsTIEElement . . . . .	132
8.2.20.1. Requested Entries . . . . .	132
8.2.21. Registry RIFT_v4/encoding/NodeTIEElement . . . . .	132
8.2.21.1. Requested Entries . . . . .	133
8.2.22. Registry RIFT_v4/encoding/PacketContent . . . . .	133
8.2.22.1. Requested Entries . . . . .	133
8.2.23. Registry RIFT_v4/encoding/PacketHeader . . . . .	133
8.2.23.1. Requested Entries . . . . .	134
8.2.24. Registry RIFT_v4/encoding/PrefixAttributes . . . . .	134
8.2.24.1. Requested Entries . . . . .	134
8.2.25. Registry RIFT_v4/encoding/PrefixTIEElement . . . . .	134
8.2.25.1. Requested Entries . . . . .	135
8.2.26. Registry RIFT_v4/encoding/ProtocolPacket . . . . .	135
8.2.26.1. Requested Entries . . . . .	135
8.2.27. Registry RIFT_v4/encoding/TIDEPacket . . . . .	135
8.2.27.1. Requested Entries . . . . .	135
8.2.28. Registry RIFT_v4/encoding/TIEElement . . . . .	135
8.2.28.1. Requested Entries . . . . .	136
8.2.29. Registry RIFT_v4/encoding/TIEHeader . . . . .	136
8.2.29.1. Requested Entries . . . . .	137
8.2.30. Registry RIFT_v4/encoding/TIEHeaderWithLifeTime . . . . .	137

8.2.30.1. Requested Entries . . . . .	137
8.2.31. Registry RIFT_v4/encoding/TIEID . . . . .	137
8.2.31.1. Requested Entries . . . . .	138
8.2.32. Registry RIFT_v4/encoding/TIEPacket . . . . .	138
8.2.32.1. Requested Entries . . . . .	138
8.2.33. Registry RIFT_v4/encoding/TIREPacket . . . . .	138
8.2.33.1. Requested Entries . . . . .	138
9. Acknowledgments . . . . .	138
10. References . . . . .	139
10.1. Normative References . . . . .	139
10.2. Informative References . . . . .	141
Appendix A. Sequence Number Binary Arithmetic . . . . .	143
Appendix B. Information Elements Schema . . . . .	144
B.1. common.thrift . . . . .	146
B.2. encoding.thrift . . . . .	152
Appendix C. Constants . . . . .	161
C.1. Configurable Protocol Constants . . . . .	161
Authors' Addresses . . . . .	163

## 1. Authors

This work is a product of a list of individuals which are all to be considered major contributors independent of the fact whether their name made it to the limited boilerplate author's list or not.

Tony Przygienda, Ed. Juniper Networks	Alankar Sharma Comcast	Pascal Thubert Cisco
Bruno Rijsman Individual	Ilya Vershkov Mellanox	Dmitry Afanasiev Yandex
Don Fedyk Individual	Alia Atlas Individual	John Drake Juniper

Table 1: RIFT Authors

## 2. Introduction

Clos [CLOS] and Fat-Tree [FATTREE] topologies have gained prominence in today's networking, primarily as result of the paradigm shift towards a centralized data-center based architecture that is poised to deliver a majority of computation and storage services in the future. Today's current routing protocols were geared towards a network with an irregular topology and low degree of connectivity originally but given they were the only available options, consequently several attempts to apply those protocols to Clos have been made. Most successfully BGP [RFC4271] [RFC7938] has been extended to this purpose, not as much due to its inherent suitability

but rather because the perceived capability to easily modify BGP and the immanent difficulties with link-state [DIJKSTRA] based protocols to optimize topology exchange and converge quickly in large scale densely meshed topologies. The incumbent protocols precondition normally extensive configuration or provisioning during bring up and re-dimensioning. This tends to be viable only for a set of organizations with according networking operation skills and budgets. For many IP fabric builders a desirable protocol would be one that auto-configures itself and deals with failures and mis-configurations with a minimum of human intervention only. Such a solution would allow local IP fabric bandwidth to be consumed in a 'standard component' fashion, i.e. provision it much faster and operate it at much lower costs than today, much like compute or storage is consumed already.

In looking at the problem through the lens of data center requirements, RIFT addresses challenges in IP fabric routing not through an incremental modification of either a link-state (distributed computation) or distance-vector (diffused computation) but rather a mixture of both, colloquially best described as "link-state towards the spine" and "distance vector towards the leaves". In other words, "bottom" levels are flooding their link-state information in the "northern" direction while each node generates under normal conditions a "default route" and floods it in the "southern" direction. This type of protocol allows naturally for highly desirable aggregation. Alas, such aggregation could blackhole traffic in cases of misconfiguration or while failures are being resolved or even cause partial network partitioning and this has to be addressed by some adequate mechanism. The approach RIFT takes is described in Section 4.2.5 and is basically based on automatic, sufficient disaggregation of prefixes in case of link and node failures.

For the visually oriented reader, Figure 1 presents a first level simplified view of the resulting information and routes on a RIFT fabric. The top of the fabric is holding in its link-state database the nodes below it and the routes to them. In the second row of the database table we indicate that partial information of other nodes in the same level is available as well. The details of how this is achieved will be postponed for the moment. When we look at the "bottom" of the fabric, the leaves, we see that the topology is basically empty and they only hold a load balanced default route to the next level under normal conditions.

The balance of this document details a dedicated IP fabric routing protocol, fills in the specification details and ultimately includes resulting security considerations.

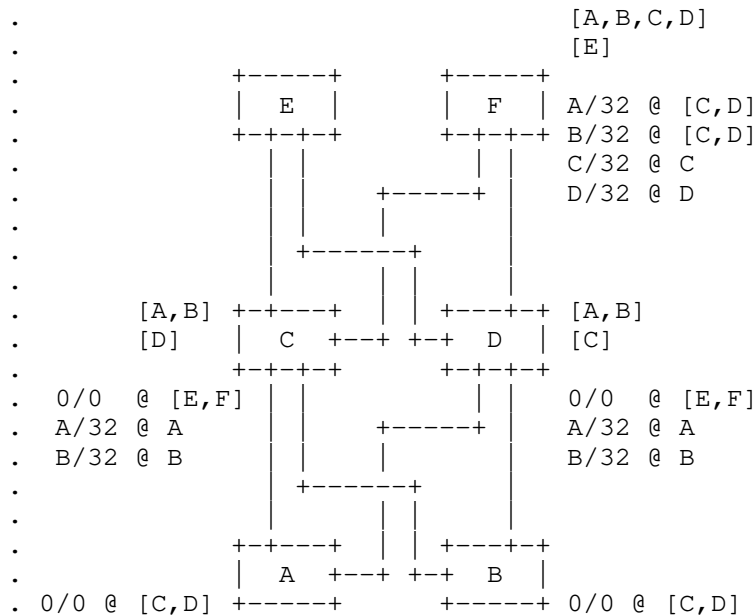


Figure 1: RIFT Information Distribution

2.1. 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 8174 [RFC8174].

3. Reference Frame

3.1. Terminology

This section presents the terminology used in this document. It is assumed that the reader is thoroughly familiar with the terms and concepts used in OSPF [RFC2328] and IS-IS [ISO10589-Second-Edition], [ISO10589] as well as the according graph theoretical concepts of shortest path first (SPF) [DIJKSTRA] computation and DAGs.

Crossbar: Physical arrangement of ports in a switching matrix without implying any further scheduling or buffering disciplines.

Clos/Fat Tree: This document uses the terms Clos and Fat Tree interchangeably whereas it always refers to a folded spine-and-leaf topology with possibly multiple Points of Delivery (PoDs) and one or multiple Top of Fabric (ToF) planes. Several modifications

such as leaf-2-leaf shortcuts and multiple level shortcuts are possible and described further in the document.

**Directed Acyclic Graph (DAG):** A finite directed graph with no directed cycles (loops). If links in Clos are considered as either being all directed towards the top or vice versa, each of such two graphs is a DAG.

**Folded Spine-and-Leaf:** In case Clos fabric input and output stages are analogous, the fabric can be "folded" to build a "superspine" or top which we will call Top of Fabric (ToF) in this document.

**Level:** Clos and Fat Tree networks are topologically partially ordered graphs and 'level' denotes the set of nodes at the same height in such a network, where the bottom level (leaf) is the level with lowest value. A node has links to nodes one level down and/or one level up. Under some circumstances, a node may have links to nodes at the same level. As footnote: Clos terminology uses often the concept of "stage" but due to the folded nature of the Fat Tree we do not use it to prevent misunderstandings.

**Superspine vs. Aggregation and Spine vs. Edge/Leaf:**

Traditional level names in 5-stages folded Clos for Level 2, 1 and 0 respectively. We normalize this language to talk about top-of-fabric (ToF), top-of-pod (ToP) and leaves.

**Zero Touch Provisioning (ZTP):** Optional RIFT mechanism which allows to derive node levels automatically based on minimum configuration (only ToF property has to be provisioned on according nodes).

**Point of Delivery (PoD):** A self-contained vertical slice or subset of a Clos or Fat Tree network containing normally only level 0 and level 1 nodes. A node in a PoD communicates with nodes in other PoDs via the Top-of-Fabric. We number PoDs to distinguish them and use PoD #0 to denote "undefined" PoD.

**Top of PoD (ToP):** The set of nodes that provide intra-PoD communication and have northbound adjacencies outside of the PoD, i.e. are at the "top" of the PoD.

**Top of Fabric (ToF):** The set of nodes that provide inter-PoD communication and have no northbound adjacencies, i.e. are at the "very top" of the fabric. ToF nodes do not belong to any PoD and are assigned "undefined" PoD value to indicate the equivalent of "any" PoD.

**Spine:** Any nodes north of leaves and south of top-of-fabric nodes. Multiple layers of spines in a PoD are possible.

**Leaf:** A node without southbound adjacencies. Its level is 0 (except cases where it is deriving its level via ZTP and is running without LEAF\_ONLY which will be explained in Section 4.2.7).

**Top-of-fabric Plane or Partition:** In large fabrics top-of-fabric switches may not have enough ports to aggregate all switches south of them and with that, the ToF is 'split' into multiple independent planes. Introduction and Section 4.1.2 explains the concept in more detail. A plane is subset of ToF nodes that see each other through south reflection or E-W links.

**Radix:** A radix of a switch is basically number of switching ports it provides. It's sometimes called fanout as well.

**North Radix:** Ports cabled northbound to higher level nodes.

**South Radix:** Ports cabled southbound to lower level nodes.

**South/Southbound and North/Northbound (Direction):**

When describing protocol elements and procedures, we will be using in different situations the directionality of the compass. I.e., 'south' or 'southbound' mean moving towards the bottom of the Clos or Fat Tree network and 'north' and 'northbound' mean moving towards the top of the Clos or Fat Tree network.

**Northbound Link:** A link to a node one level up or in other words, one level further north.

**Southbound Link:** A link to a node one level down or in other words, one level further south.

**East-West Link:** A link between two nodes at the same level. East-West links are normally not part of Clos or "fat-tree" topologies.

**Leaf shortcuts (L2L):** East-West links at leaf level will need to be differentiated from East-West links at other levels.

**Routing on the host (RotH):** Modern data center architecture variant where servers/leaves are multi-homed and consecutively participate in routing.

**Northbound representation:** Subset of topology information flooded towards higher levels of the fabric.

**Southbound representation:** Subset of topology information sent towards a lower level.

**South Reflection:** Often abbreviated just as "reflection" it defines a mechanism where South Node TIEs are "reflected" from the level south back up north to allow nodes in the same level without E-W links to "see" each other's node TIEs.

**TIE:** This is an acronym for a "Topology Information Element". TIEs are exchanged between RIFT nodes to describe parts of a network such as links and address prefixes, in a fashion similar to ISIS LSPs or OSPF LSAs. A TIE has always a direction and a type. We will talk about North TIEs (sometimes abbreviated as N-TIEs) when talking about TIEs in the northbound representation and South-TIEs (sometimes abbreviated as S-TIEs) for the southbound equivalent. TIEs have different types such as node and prefix TIEs.

**Node TIE:** This stands as acronym for a "Node Topology Information Element" that contains all adjacencies the node discovered and information about node itself. Node TIE should NOT be confused with a North TIE since "node" defines the type of TIE rather than its direction.

**Prefix TIE:** This is an acronym for a "Prefix Topology Information Element" and it contains all prefixes directly attached to this node in case of a North TIE and in case of South TIE the necessary default routes the node advertises southbound.

**Key Value TIE:** A South TIE that is carrying a set of key value pairs [DYNAMO]. It can be used to distribute information in the southbound direction within the protocol.

**TIDE:** Topology Information Description Element, equivalent to CSNP in ISIS.

**TIRE:** Topology Information Request Element, equivalent to PSNP in ISIS. It can both confirm received and request missing TIEs.

**De-aggregation/Disaggregation:** Process in which a node decides to advertise more specific prefixes Southwards, either positively to attract the corresponding traffic, or negatively to repel it. Disaggregation is performed to prevent black-holing and suboptimal routing to the more specific prefixes.

**LIE:** This is an acronym for a "Link Information Element", largely equivalent to HELLOs in IGPs and exchanged over all the links between systems running RIFT to form three way adjacencies.

**Flood Repeater (FR):** A node can designate one or more northbound neighbor nodes to be flood repeaters. The flood repeaters are responsible for flooding northbound TIEs further north. They are



similar to MPR in OSLR. The document sometimes calls them flood leaders as well.

**Bandwidth Adjusted Distance (BAD):** Each RIFT node can calculate the amount of northbound bandwidth available towards a node compared to other nodes at the same level and can modify the route distance accordingly to allow for the lower level to adjust their load balancing towards spines.

**Overloaded:** Applies to a node advertising 'overload' attribute as set. The semantics closely follow the meaning of the same attribute in [ISO10589-Second-Edition].

**Interface:** A layer 3 entity over which RIFT control packets are exchanged.

**Three-Way Adjacency:** RIFT tries to form a unique adjacency over an interface and exchange local configuration and necessary ZTP information. An adjacency is only advertised in node TIEs and used for computations after it achieved three-way state, i.e. both routers reflected each other in LIEs including relevant security information. LIEs before three-way state is reached may carry ZTP related information already.

**Bi-directional Adjacency:** Bidirectional adjacency is an adjacency where nodes of both sides of the adjacency advertised it in the node TIEs with the correct levels and system IDs. Bi-directionality is used to check in different algorithms whether the link should be included.

**Neighbor:** Once a three-way adjacency has been formed a neighborhood relationship contains the neighbor's properties. Multiple adjacencies can be formed to a remote node via parallel interfaces but such adjacencies are NOT sharing a neighbor structure. Saying "neighbor" is thus equivalent to saying "a three-way adjacency".

**Cost:** The term signifies the weighted distance between two neighbors.

**Distance:** Sum of costs (bound by infinite distance) between two nodes.

**Shortest-Path First (SPF):** A well-known graph algorithm attributed to Dijkstra that establishes a tree of shortest paths from a source to destinations on the graph. We use SPF acronym due to its familiarity as general term for the node reachability calculations RIFT can employ to ultimately calculate routes of which Dijkstra algorithm is one.

North SPF (N-SPF): A reachability calculation that is progressing northbound, as example SPF that is using South Node TIEs only. Normally it progresses a single hop only and installs default routes.

South SPF (S-SPF): A reachability calculation that is progressing southbound, as example SPF that is using North Node TIEs only.

Security Envelope RIFT packets are flooded within an authenticated security envelope that allows to protect the integrity of information a node accepts.

### 3.2. Topology

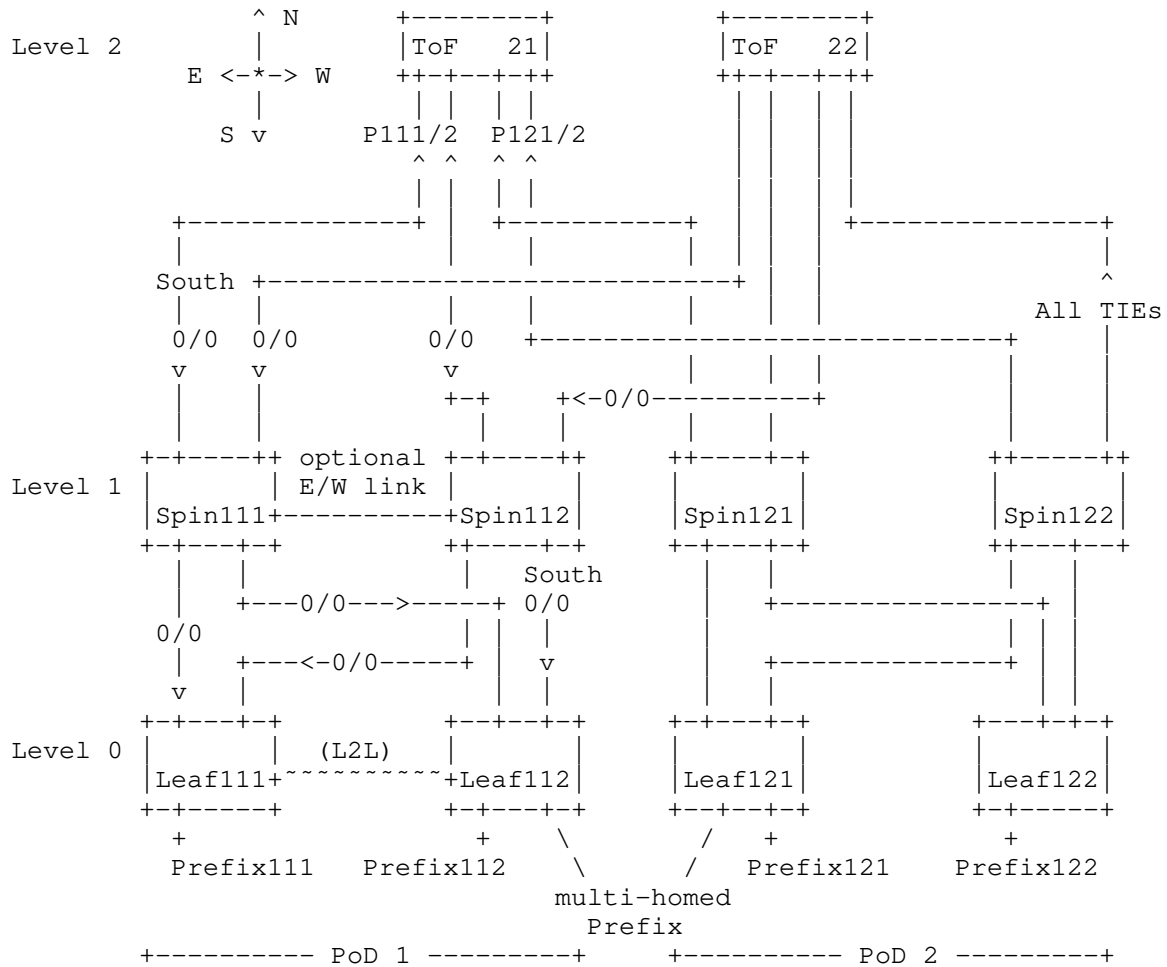


Figure 2: A Three Level Spine-and-Leaf Topology

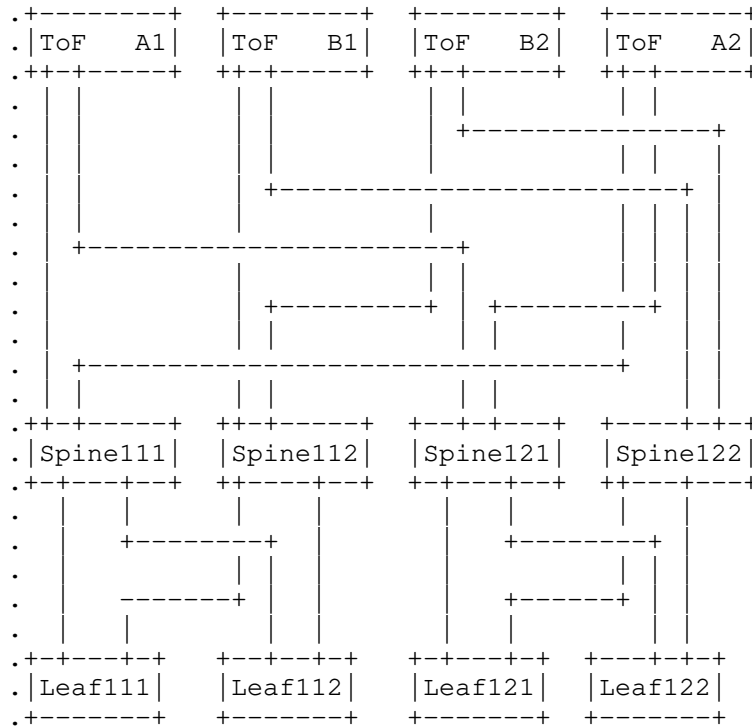


Figure 3: Topology with Multiple Planes

We will use topology in Figure 2 (called commonly a fat tree/network in modern IP fabric considerations [VAHDAT08] as homonym to the original definition of the term [FATREE]) in all further considerations. This figure depicts a generic "single plane fat-tree" and the concepts explained using three levels apply by induction to further levels and higher degrees of connectivity. Further, this document will deal also with designs that provide only sparser connectivity and "partitioned spines" as shown in Figure 3 and explained further in Section 4.1.2.

4. RIFT: Routing in Fat Trees

We present here a detailed outline of a protocol optimized for Routing in Fat Trees (RIFT) that in most abstract terms has many properties of a modified link-state protocol [RFC2328][ISO10589-Second-Edition] when distributing information northbound and distance vector [RFC4271] protocol when distributing information southbound. While this is an unusual combination, it does quite naturally exhibit the desirable properties we seek.

## 4.1. Overview

### 4.1.1. Properties

The most singular property of RIFT is that it floods flat link-state information northbound only so that each level obtains the full topology of levels south of it. Link-State information is, with some exceptions, never flooded East-West or back South again. Exceptions like south reflection is explained in detail in Section 4.2.5.1 and east-west flooding at ToF level in multi-plane fabrics is outlined in Section 4.1.2. In southbound direction, the protocol operates like a "fully summarizing, unidirectional" path vector protocol or rather a distance vector with implicit split horizon. Routing information, normally just the default route, propagates one hop south and is 're-advertised' by nodes at next lower level. However, RIFT uses flooding in the southern direction as well to avoid the overhead of building an update per adjacency. We omit describing the East-West direction for the moment.

Those information flow constraints create not only an anisotropic protocol (i.e. the information is not distributed "evenly" or "clumped" but summarized along the N-S gradient) but also a "smooth" information propagation where nodes do not receive the same information from multiple directions at the same time. Normally, accepting the same reachability on any link, without understanding its topological significance, forces tie-breaking on some kind of distance metric. And such tie-breaking leads ultimately in hop-by-hop forwarding to shortest paths only. In contrast to that, RIFT, under normal conditions, does not need to tie-break same reachability information from multiple directions. Its computation principles (south forwarding direction is always preferred) leads to valley-free forwarding behavior. And since valley free routing is loop-free, it can use all feasible paths which is another highly desirable property if available bandwidth should be utilized to the maximum extent possible.

To account for the "northern" and the "southern" information split the link state database is partitioned accordingly into "north representation" and "south representation" TIEs. In simplest terms the North TIEs contain a link state topology description of lower levels and and South TIEs carry simply default routes towards the level above. This oversimplified view will be refined gradually in following sections while introducing protocol procedures and state machines at the same time.

#### 4.1.2. Generalized Topology View

This section will shed some light on the topologies RIFT addresses, including multi plane fabrics and their implications. Readers that are only interested in single plane designs, i.e. all top-of-fabric nodes being topologically equal and initially connected to all the switches at the level below them, can skip the rest of Section 4.1.2 and resulting Section 4.2.5.2 as well.

It is quite difficult to visualize multi plane design, which are effectively multi-dimensional switching matrices. To cope with that, we will introduce a methodology allowing us to depict the connectivity in two-dimensional pictures. Further, we will leverage the fact that we are dealing basically with stacked crossbar fabrics where ports align "on top of each other" in a regular fashion.

A word of caution to the reader; at this point it should be observed that the language used to describe Clos variations, especially in multi-plane designs, varies widely between sources. This description follows the terminology introduced in Section 3.1. It is unavoidable to have it present to be able to follow the rest of this section correctly.

##### 4.1.2.1. Terminology

This section describes the terminology and acronyms used in the rest of the text.

P: Denotes the number of PoDs in a topology.

S: Denotes the number of ToF nodes in a topology.

K: Denotes the number of ports in radix of a switch pointing north or south. Further,  $K_{LEAF}$  denotes number of ports pointing south, i.e. towards leaves, and  $K_{TOP}$  for number of ports pointing north towards a higher spine level. To simplify the visual aids, notations and further considerations, K will be mostly set to  $Radix/2$ .

ToF Plane: Set of ToFs that are aware of each other by means of south reflection. We number planes by capital letters, e.g. plane A.

N: Denote the number of independent ToF planes in a topology.

R: Denotes a redundancy factor, i.e. number of connections a spine has towards a ToF plane. In single plane design  $K_{TOP}$  is equal to R.

Fallen Leaf: A fallen leaf in a plane Z is a switch that lost all connectivity northbound to Z.

#### 4.1.2.2. Clos as Crossed Crossbars

The typical topology for which RIFT is defined is built of P number of PoDs and connected together by S number of ToF nodes. A PoD node has K number of ports (also called Radix). We consider half of them ( $K=Radix/2$ ) as connecting host devices from the south, and the other half connecting to interleaved PoD Top-Level switches to the north. Ratio K can be chosen differently without loss of generality when port speeds differ or the fabric is oversubscribed but  $K=R/2$  allows for more readable representation whereby there are as many ports facing north as south on any intermediate node. We represent a node hence in a schematic fashion with ports "sticking out" to its north and south rather than by the usual real-world front faceplate designs of the day.

Figure 4 provides a view of a leaf node as seen from the north, i.e. showing ports that connect northbound. For lack of a better symbol, we have chosen to use the "o" as ASCII visualisation of a single port. In this example, K\_LEAF has 6 ports. Observe that the number of PoDs is not related to Radix unless the ToF Nodes are constrained to be the same as the PoD nodes in a particular deployment.

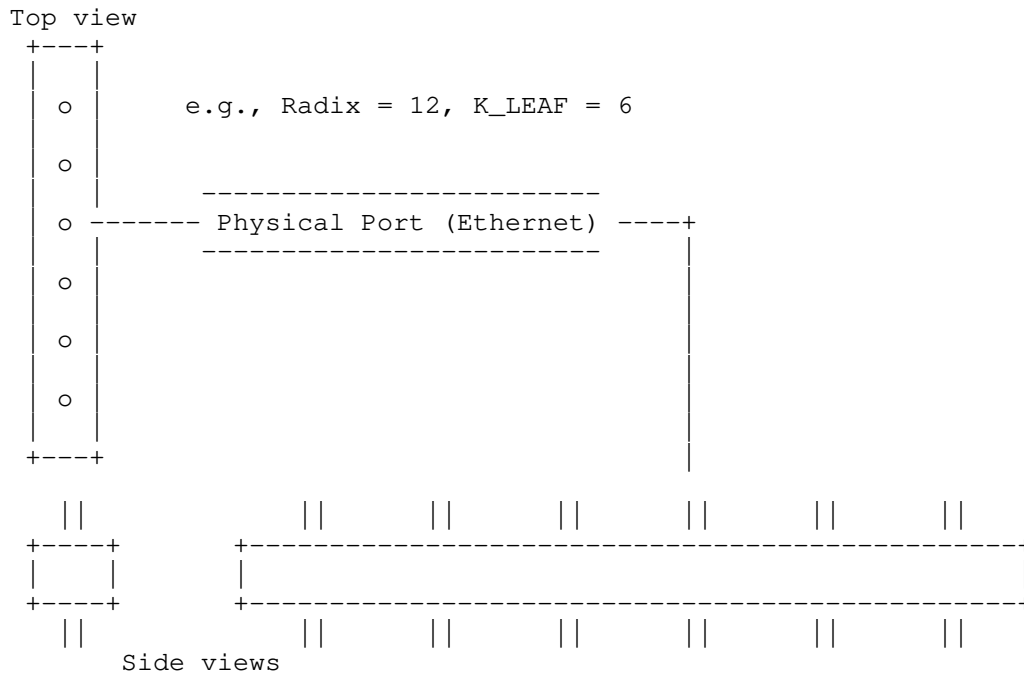


Figure 4: A Leaf Node,  $K_{LEAF}=6$

The Radix of a PoD's top node may be different than that of the leaf node. Though, more often than not, a same type of node is used for both, effectively forming a square ( $K*K$ ). In general case, we could have switches with  $K_{TOP}$  southern ports on nodes at the top of the PoD which are not necessarily the same as  $K_{LEAF}$ . For instance, in the representations below, we pick a 6 port  $K_{LEAF}$  and a 8 port  $K_{TOP}$ . In order to form a crossbar, we need  $K_{TOP}$  Leaf Nodes as illustrated in Figure 5.



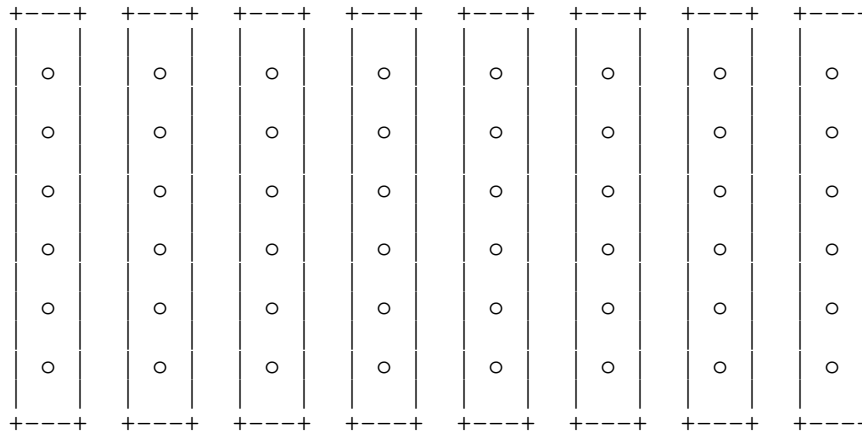


Figure 5: Southern View of a PoD, K\_TOP=8

As further visualized in Figure 6 the K\_TOP Leaf Nodes are fully interconnected with the K\_LEAF PoD-top nodes, providing connectivity that can be represented as a crossbar when "looked at" from the north. The result is that, in the absence of a failure, a packet entering the PoD from the north on any port can be routed to any port in the south of the PoD and vice versa. And that is precisely why it makes sense to talk about a "switching matrix".

E<-\*->W

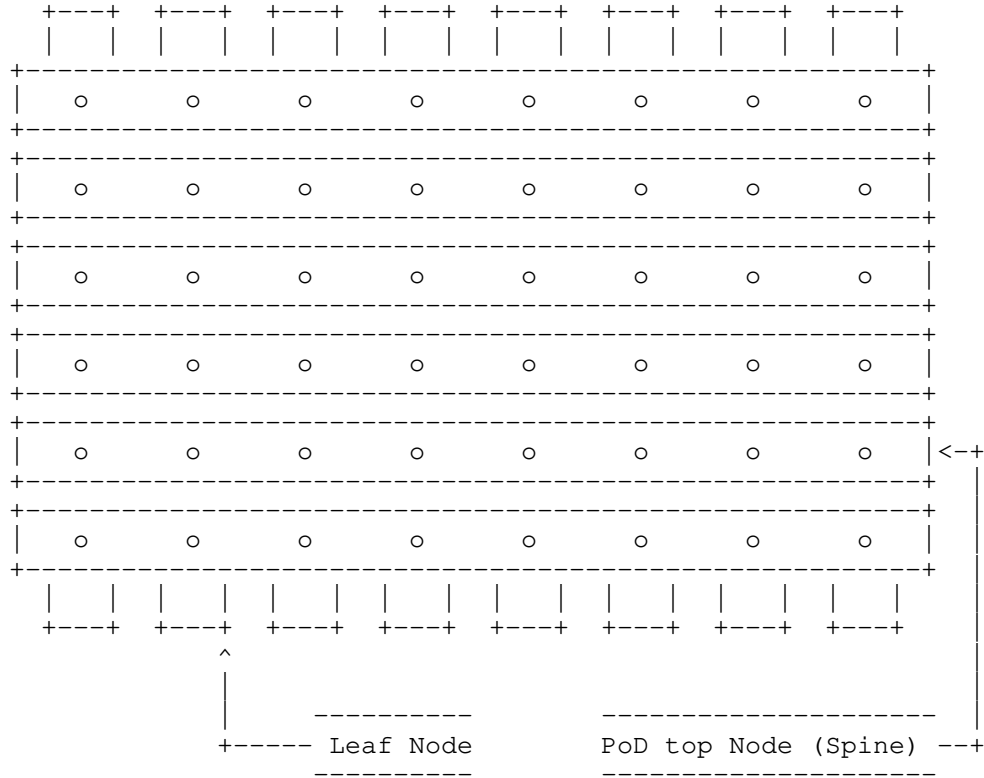


Figure 6: Northern View of a PoD's Spines, K\_TOP=8

Side views of this PoD is illustrated in Figure 7 and Figure 8.

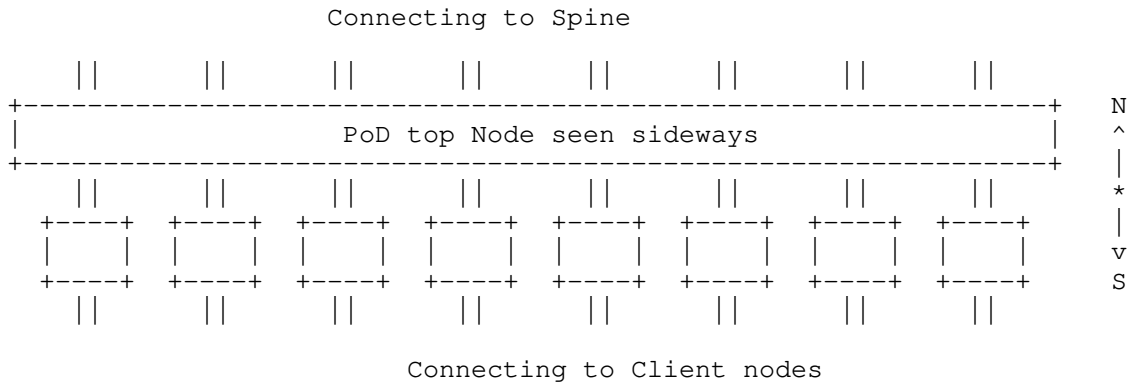


Figure 7: Side View of a PoD,  $K_{TOP}=8$ ,  $K_{LEAF}=6$

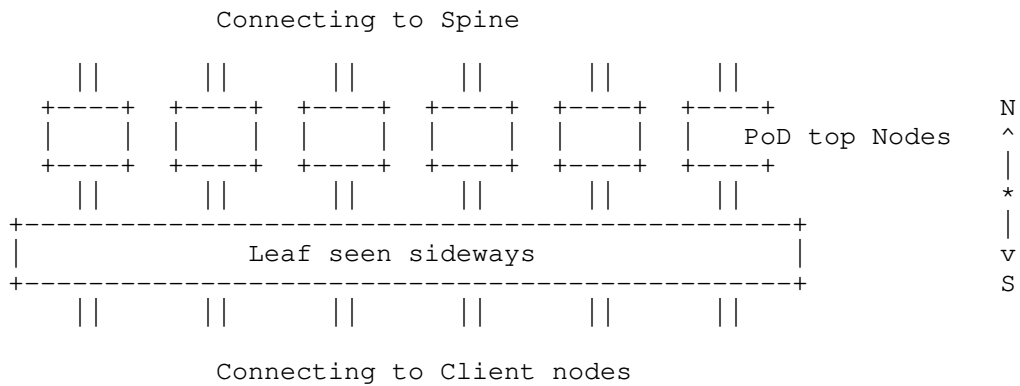


Figure 8: Other Side View of a PoD,  $K_{TOP}=8$ ,  $K_{LEAF}=6$ , 90o turn in E-W Plane

As next step, let us observe that a resulting PoD can be abstracted as a bigger node with a number K of  $K_{POD}= K_{TOP} * K_{LEAF}$ , and the design can recurse.

It will be critical at this point that, before progressing further, the concept and the picture of "crossed crossbars" is clear. Else, the following considerations might be difficult to comprehend.

To continue, the PoDs are interconnected with each other through a Top-of-Fabric (ToF) node at the very top or the north edge of the fabric. The resulting ToF is NOT partitioned if, and only if (IIF), every PoD top level node (spine) is connected to every ToF Node.

This topology is also referred to as a single plane configuration and is quite popular due to its simplicity. In order to reach a 1:1 connectivity ratio between the ToF and the leaves, it results that there are  $K_{TOP}$  ToF nodes, because each port of a ToP node connects to a different ToF node, and  $K_{LEAF}$  ToP nodes for the same reason. Consequently, it will take  $(P * K_{LEAF})$  ports on a ToF node to connect to each of the  $K_{LEAF}$  ToP nodes of the  $P$  PoDs, as shown in Figure 9.

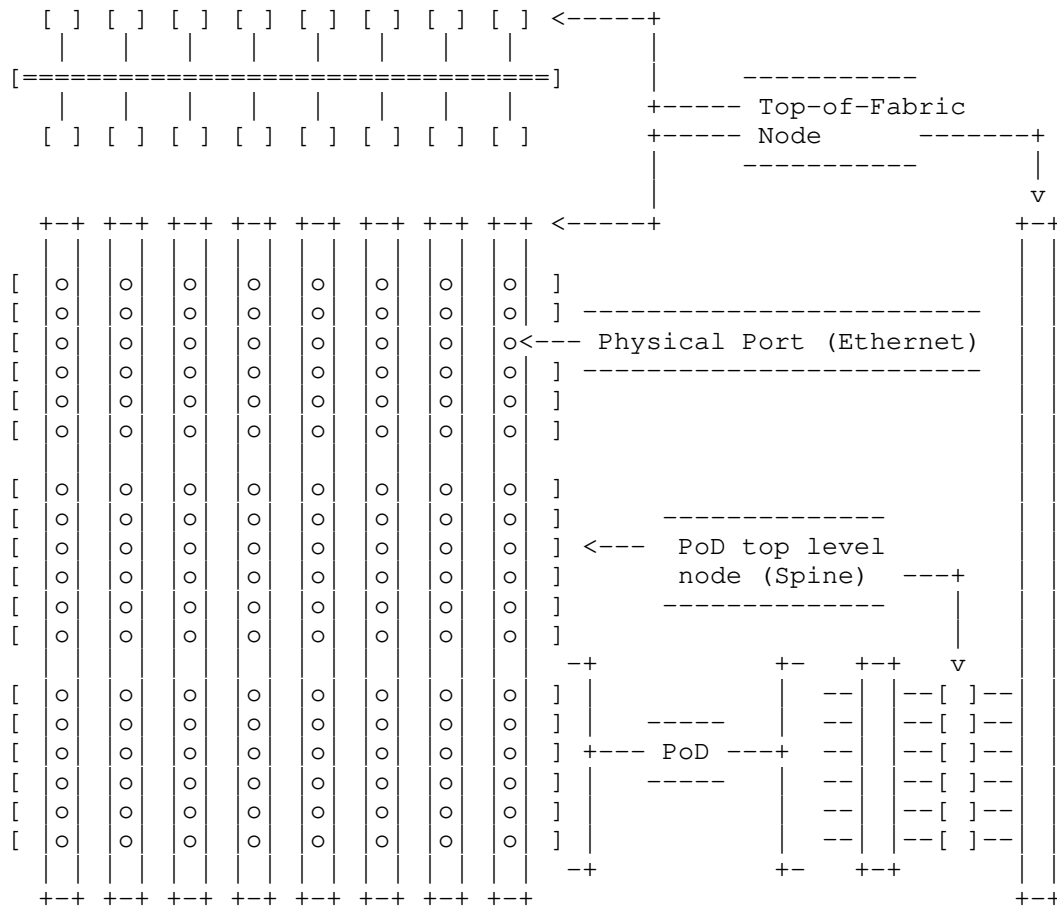


Figure 9: Fabric Spines and TOFs in Single Plane Design, 3 PoDs

The top view can be collapsed into a third dimension where the hidden depth index is representing the PoD number. We can then show one PoD as a class of PoDs and hence save one dimension in our

representation. The Spine Node expands in the depth and the vertical dimensions, whereas the PoD top level Nodes are constrained, in horizontal dimension. A port in the 2-D representation represents effectively the class of all the ports at the same position in all the PoDs that are projected in its position along the depth axis. This is shown in Figure 10.

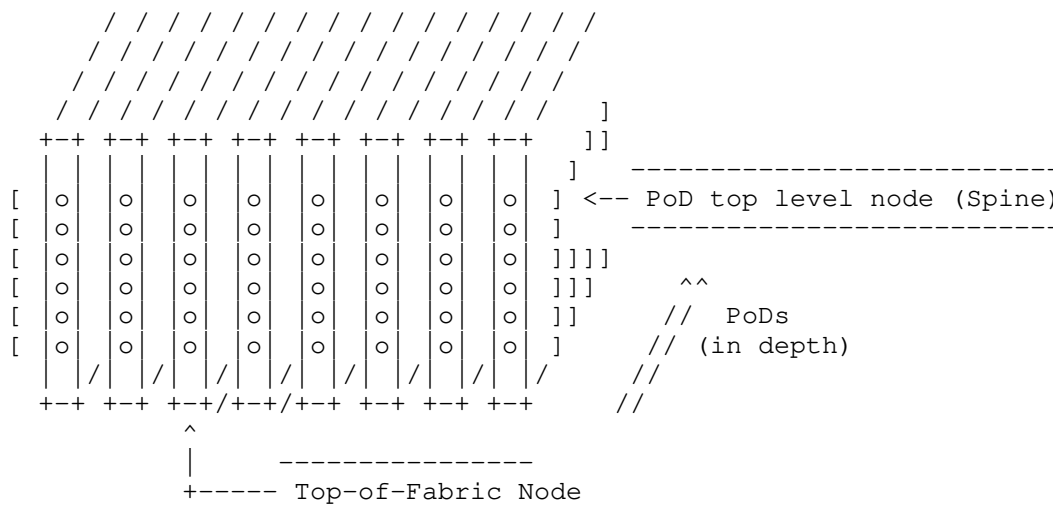


Figure 10: Collapsed Northern View of a Fabric for Any Number of PoDs

As simple as single plane deployment is it introduces a limit due to the bound on the available radix of the ToF nodes that has to be at least  $P * K\_LEAF$ . Nevertheless, we will see that a distinct advantage of a connected or non-partitioned Top-of-Fabric is that all failures can be resolved by simple, non-transitive, positive disaggregation (i.e. nodes advertising more specific prefixes with the default to the level below them that is however not propagated further down the fabric) as described in Section 4.2.5.1 . In other words; non-partitioned ToF nodes can always reach nodes below or withdraw the routes from PoDs they cannot reach unambiguously. And with this, positive disaggregation can heal all failures and still allow all the ToF nodes to see each other via south reflection. Disaggregation will be explained in further detail in Section 4.2.5.

In order to scale beyond the "single plane limit", the Top-of-Fabric can be partitioned by a N number of identically wired planes where N is an integer divider of  $K\_LEAF$ . The 1:1 ratio and the desired symmetry are still served, this time with  $(K\_TOP * N)$  ToF nodes, each of  $(P * K\_LEAF / N)$  ports.  $N=1$  represents a non-partitioned Spine and  $N=K\_LEAF$  is a maximally partitioned Spine. Further, if R is any

integer divisor of  $K\_LEAF$ , then  $N=K\_LEAF/R$  is a feasible number of planes and  $R$  a redundancy factor. It proves convenient for deployments to use a radix for the leaf nodes that is a power of 2 so they can pick a number of planes that is a lower power of 2. The example in Figure 11 splits the Spine in 2 planes with a redundancy factor  $R=3$ , meaning that there are 3 non-intersecting paths between any leaf node and any ToF node. A ToF node must have, in this case, at least  $3 \cdot P$  ports, and be directly connected to 3 of the 6 PoD-ToP nodes (spines) in each PoD.

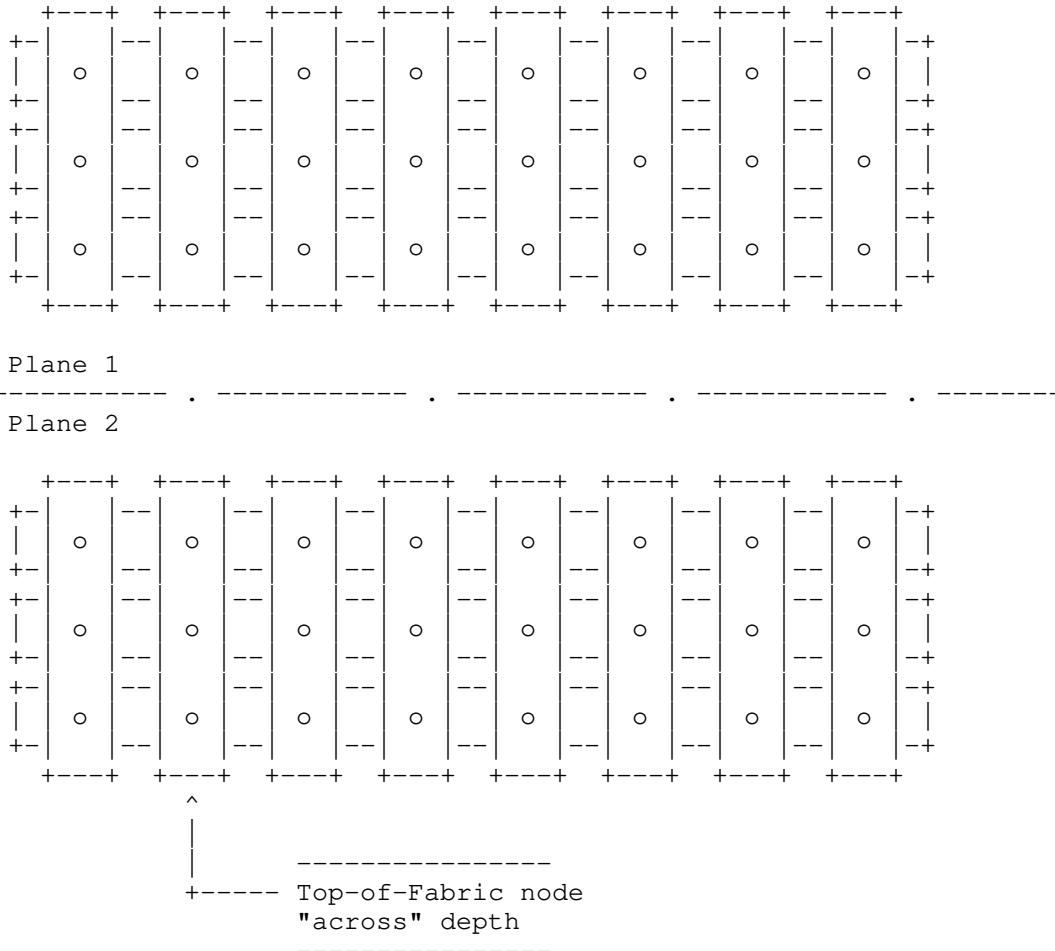


Figure 11: Northern View of a Multi-Plane ToF Level,  $K\_LEAF=6$ ,  $N=2$

At the extreme end of the spectrum it is even possible to fully partition the spine with  $N = K\_LEAF$  and  $R=1$ , while maintaining

connectivity between each leaf node and each Top-of-Fabric node. In that case the ToF node connects to a single Port per PoD, so it appears as a single port in the projected view represented in Figure 12. The number of ports required on the Spine Node is more or equal to  $P$ , the number of PoDs.

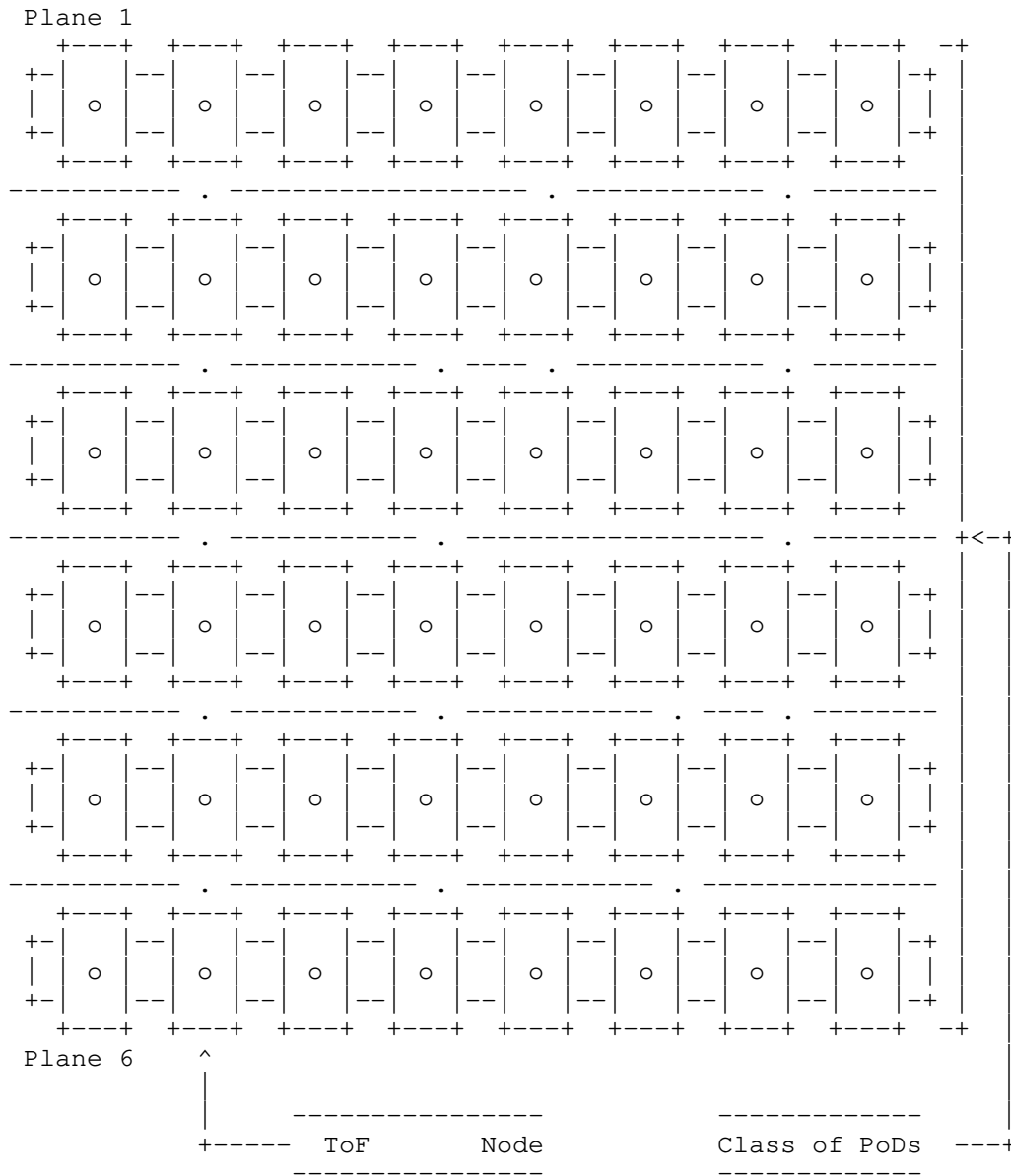


Figure 12: Northern View of a Maximally Partitioned ToF Level, R=1



#### 4.1.3. Fallen Leaf Problem

As mentioned earlier, RIFT exhibits an anisotropic behavior tailored for fabrics with a North / South orientation and a high level of interleaving paths. A non-partitioned fabric makes a total loss of connectivity between a Top-of-Fabric node at the north and a leaf node at the south a very rare but yet possible occasion that is fully healed by positive disaggregation as described in Section 4.2.5.1. In large fabrics or fabrics built from switches with low radix, the ToF ends often being partitioned in planes which makes the occurrence of having a given leaf being only reachable from a subset of the ToF nodes more likely to happen. This makes some further considerations necessary.

We define a "Fallen Leaf" as a leaf that can be reached by only a subset, but not all, of Top-of-Fabric nodes due to missing connectivity. If  $R$  is the redundancy factor, then it takes at least  $R$  breakages to reach a "Fallen Leaf" situation.

In a maximally partitioned fabric, the redundancy factor is  $R=1$ , so any breakage in the fabric may cause one or more fallen leaves. However, not all cases require disaggregation. The following cases do not require particular action in such scenario:

If a southern link on a node goes down, then connectivity through that node is lost for all nodes south of it. There is no need to disaggregate since the connectivity to this node is lost for all spine nodes in a same fashion.

If a ToF Node goes down, then northern traffic towards it is routed via alternate ToF nodes in the same plane and there is no need to disaggregate routes.

In a general manner, the mechanism of non-transitive positive disaggregation is sufficient when the disaggregating ToF nodes collectively connect to all the ToP nodes in the broken plane. This happens in the following case:

If the breakage is the last northern link from a ToP node to a ToF node going down, then the fallen leaf problem affects only The ToF node, and the connectivity to all the nodes in the PoD is lost from that ToF node. This can be observed by other ToF nodes within the plane where the ToP node is located and positively disaggregated within that plane.

On the other hand, there is a need to disaggregate the routes to Fallen Leaves in a transitive fashion, all the way to the other leaves in the following cases:

- o If the breakage is the last northern link from a leaf node within a plane (there is only one such link in a maximally partitioned fabric) that goes down, then connectivity to all unicast prefixes attached to the leaf node is lost within the plane where the link is located. Southern Reflection by a leaf node, e.g., between ToP nodes, if the PoD has only 2 levels, happens in between planes, allowing the ToP nodes to detect the problem within the PoD where it occurs and positively disaggregate. The breakage can be observed by the ToF nodes in the same plane through the North flooding of TIEs from the ToP nodes. The ToF nodes however need to be aware of all the affected prefixes for the negative, possibly transitive disaggregation to be fully effective (i.e. a node advertising in control plane that it cannot reach a certain more specific prefix than default whereas such disaggregation must in extreme condition propagate further down southbound). The problem can also be observed by the ToF nodes in the other planes through the flooding of North TIEs from the affected leaf nodes, together with non-node North TIEs which indicate the affected prefixes. To be effective in that case, the positive disaggregation must reach down to the nodes that make the plane selection, which are typically the ingress leaf nodes. The information is not useful for routing in the intermediate levels.
- o If the breakage is a ToP node in a maximally partitioned fabric - in which case it is the only ToP node serving the plane in that PoD - goes down, then the connectivity to all the nodes in the PoD is lost within the plane where the ToP node is located. Consequently, all leaves of the PoD fall in this plane. Since the Southern Reflection between the ToF nodes happens only within a plane, ToF nodes in other planes cannot discover fallen leaves in a different plane. They also cannot determine beyond their local plane whether a leaf node that was initially reachable has become unreachable. As the breakage can be observed by the ToF nodes in the plane where the breakage happened, the ToF nodes in the plane need to be aware of all the affected prefixes for the negative disaggregation to be fully effective. The problem can also be observed by the ToF nodes in the other planes through the flooding of North TIEs from the affected leaf nodes, if there are only 3 levels and the ToP nodes are directly connected to the leaf nodes, and then again it can only be effective if it is propagated transitively to the leaf, and useless above that level.

For the sake of easy comprehension let us roll the abstractions back into a simple example and observe that in Figure 3 the loss of link Spine 122 to Leaf 122 will make Leaf 122 a fallen leaf for Top-of-Fabric plane B. Worse, if the cabling was never present in first place, plane B will not even be able to know that such a fallen leaf

exists. Hence partitioning without further treatment results in two grave problems:

- o Leaf 111 trying to route to Leaf 122 MUST choose Spine 111 in plane A as its next hop since plane B will inevitably blackhole the packet when forwarding using default routes or do excessive bow tying. This information must be in its routing table.
- o Any kind of "flooding" or distance vector trying to deal with the problem by distributing host routes will be able to converge only using paths through leaves. The flooding of information on Leaf 122 would have to go up to Top-of-Fabric A and then "loopback" over other leaves to ToF B leading in extreme cases to traffic for Leaf 122 when presented to plane B taking an "inverted fabric" path where leaves start to serve as TOFs, at least for the duration of a protocol's convergence.

#### 4.1.4. Discovering Fallen Leaves

As illustrated later, and without further proof, the way to deal with fallen leaves in multi-plane designs, when aggregation is used, is that RIFT requires all the ToF nodes to share the same north topology database. This happens naturally in single plane design by the means of northbound flooding and south reflection but needs additional considerations in multi-plane fabrics. To satisfy this RIFT, in multi-plane designs, relies at the ToF level on ring interconnection of switches in multiple planes. Other solutions are possible but they either need more cabling or end up having much longer flooding paths and/or single points of failure.

In detail, by reserving two ports on each Top-of-Fabric node it is possible to connect them together by interplane bi-directional rings as illustrated in Figure 13. The rings will be used to exchange full north topology information between planes. All ToFs having same north topology allows by the means of transitive, negative disaggregation described in Section 4.2.5.2 to efficiently fix any possible fallen leaf scenario. Somewhat as a side-effect, the exchange of information fulfills the ask to present full view of the fabric topology at the Top-of-Fabric level, without the need to collate it from multiple points by additional complexity of technologies like [RFC7752].

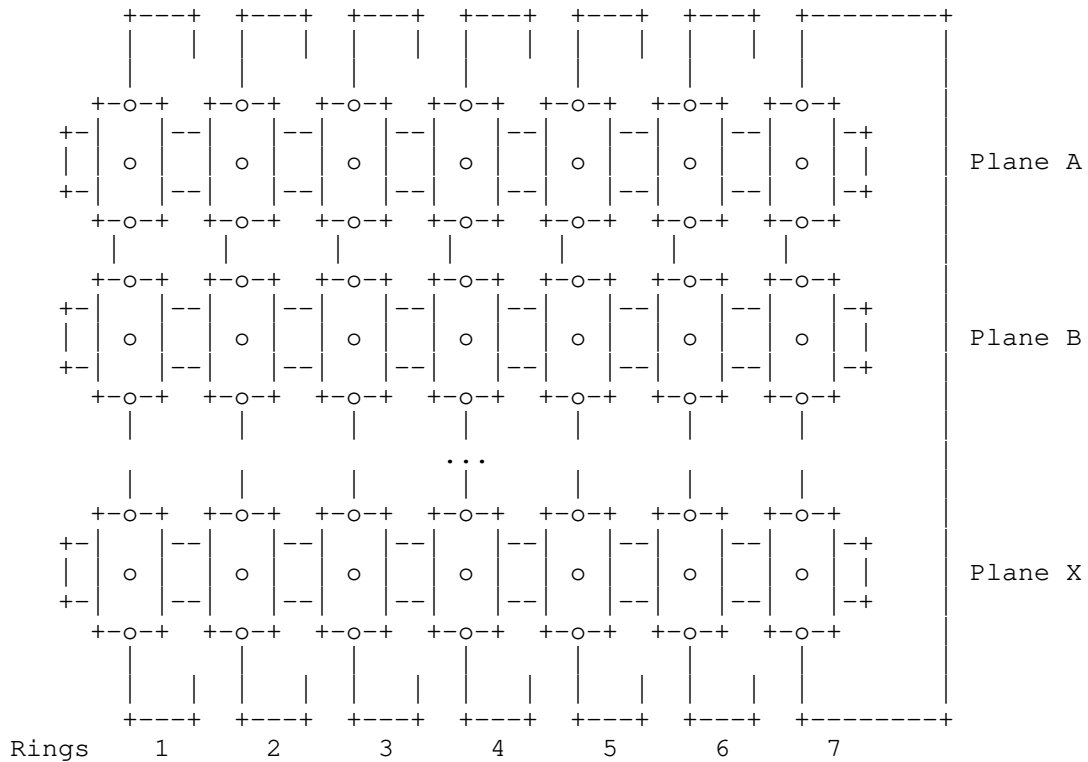


Figure 13: Connecting Top-of-Fabric Nodes Across Planes by Rings

4.1.5. Addressing the Fallen Leaves Problem

One consequence of the "Fallen Leaf" problem is that some prefixes attached to the fallen leaf become unreachable from some of the ToF nodes. RIFT proposes two methods to address this issue, the positive and the negative disaggregation. Both methods flood South TIEs to advertise the impacted prefix(es).

When used for the operation of disaggregation, a positive South TIE, as usual, indicates reachability to a prefix of given length and all addresses subsumed by it. In contrast, a negative route advertisement indicates that the origin cannot route to the advertised prefix.

The positive disaggregation is originated by a router that can still reach the advertised prefix, and the operation is not transitive. In other words, the receiver does not generate its own flooding south as a consequence of receiving positive disaggregation advertisements

from a higher level node. The effect of a positive disaggregation is that the traffic to the impacted prefix will follow the longest match and will be limited to the northbound routers that advertised the more specific route.

In contrast, the negative disaggregation can be transitive, and is propagated south when all the possible routes have been advertised as negative exceptions. A negative route advertisement is only actionable when the negative prefix is aggregated by a positive route advertisement for a shorter prefix. In such case, the negative advertisement "punches out a hole" in the positive route in the routing table, making the positive prefix reachable through the originator with the special consideration of the negative prefix removing certain next hop neighbors.

When the ToF is not partitioned, the collective southern flooding of the positive disaggregation by the ToF nodes that can still reach the impacted prefix is in general enough to cover all the switches at the next level south, typically the ToP nodes. If all those switches are aware of the disaggregation, they collectively create a ceiling that intercepts all the traffic north and forwards it to the ToF nodes that advertised the more specific route. In that case, the positive disaggregation alone is sufficient to solve the fallen leaf problem.

On the other hand, when the fabric is partitioned in planes, the positive disaggregation from ToF nodes in different planes do not reach the ToP switches in the affected plane and cannot solve the fallen leaves problem. In other words, a breakage in a plane can only be solved in that plane. Also, the selection of the plane for a packet typically occurs at the leaf level and the disaggregation must be transitive and reach all the leaves. In that case, the negative disaggregation is necessary. The details on the RIFT approach to deal with fallen leaves in an optimal way are specified in Section 4.2.5.2.

#### 4.2. Specification

This section specifies the protocol in a normative fashion by either prescriptive procedures or behavior defined by Finite State Machines (FSM).

Some FSM figures are provided as [DOT] description due to limitations of ASCII art.

"On Entry" actions on FSM state are performed every time and right before the according state is entered, i.e. after any transitions from previous state.

"On Exit" actions are performed every time and immediately when a state is exited, i.e. before any transitions towards target state are performed.

Any attempt to transition from a state towards another on reception of an event where no action is specified MUST be considered an unrecoverable error.

The FSMs and procedures are normative in the sense that an implementation MUST implement them either literally or an implementation MUST exhibit externally observable behavior that is identical to the execution of the specified FSMs.

Where a FSM representation is inconvenient, i.e. the amount of procedures and kept state exceeds the amount of transitions, we defer to a more procedural description on data structures.

#### 4.2.1. Transport

All packet formats are defined in Thrift [thrift] models in Appendix B.

The serialized model is carried in an envelope within a UDP frame that provides security and allows validation/modification of several important fields without de-serialization for performance and security reasons.

#### 4.2.2. Link (Neighbor) Discovery (LIE Exchange)

RIFT LIE exchange auto-discovers neighbors, negotiates ZTP parameters and discovers miscablings. It uses a three-way handshake mechanism which is a cleaned up version of [RFC5303]. Observe that for easier comprehension the terminology of one/two and three-way states does NOT align with OSPF or ISIS FSMs albeit they use roughly same mechanisms. The formation progresses under normal conditions from one-way to two-way and then three-way state at which point it is ready to exchange TIEs per Section 4.2.3.

LIE exchange happens over well-known administratively locally scoped and configured or otherwise well-known IPv4 multicast address [RFC2365] and/or link-local multicast scope [RFC4291] for IPv6 [RFC8200] using a configured or otherwise a well-known destination UDP port defined in Appendix C.1. LIEs SHOULD be sent with an IPv4 Time to Live (TTL) / IPv6 Hop Limit (HL) of 1 to prevent RIFT information reaching beyond a single L3 next-hop in the topology. LIEs SHOULD be sent with network control precedence.

Originating port of the LIE has no further significance other than identifying the origination point. LIEs are exchanged over all links running RIFT.

An implementation MAY listen and send LIEs on IPv4 and/or IPv6 multicast addresses. A node MUST NOT originate LIEs on an address family if it does not process received LIEs on that family. LIEs on same link are considered part of the same negotiation independent of the address family they arrive on. Observe further that the LIE source address may not identify the peer uniquely in unnumbered or link-local address cases so the response transmission MUST occur over the same interface the LIEs have been received on. A node MAY use any of the adjacency's source addresses it saw in LIEs on the specific interface during adjacency formation to send LIEs. That implies that an implementation MUST be ready to accept LIEs on all addresses it used as source of LIE frames.

A three-way adjacency over any address family implies support for IPv4 forwarding if the `'v4_forwarding_capable'` flag is set to true and a node can use [RFC5549] type of forwarding in such a situation. It is expected that the whole fabric supports the same type of forwarding of address families on all the links. Operation of a fabric where only some of the links are supporting forwarding on an address family and others do not is outside the scope of this specification.

The protocol does NOT support selective disabling of address families, disabling v4 forwarding capability or any local address changes in three-way state, i.e. if a link has entered three-way IPv4 and/or IPv6 with a neighbor on an adjacency and it wants to stop supporting one of the families or change any of its local addresses or stop v4 forwarding, it has to tear down and rebuild the adjacency. It also has to remove any information it stored about the adjacency such as LIE source addresses seen.

Unless ZTP as described in Section 4.2.7 is used, each node is provisioned with the level at which it is operating. It MAY be also provisioned with its PoD. If any of those values is undefined, then accordingly a default level and/or an "undefined" PoD are assumed. This means that leaves do not need to be configured at all if initial configuration values are all left at "undefined" value. Nodes above ToP MUST remain at "any" PoD value which has the same value as "undefined" PoD. This information is propagated in the LIEs exchanged.

Further definitions of leaf flags are found in Section 4.2.7 given they have implications in terms of level and adjacency forming here.

A node tries to form a three-way adjacency if and only if

1. the node is in the same PoD or either the node or the neighbor advertises "undefined/any" PoD membership (PoD# = 0) AND
2. the neighboring node is running the same MAJOR schema version AND
3. the neighbor is not member of some PoD while the node has a northbound adjacency already joining another PoD AND
4. the neighboring node uses a valid System ID AND
5. the neighboring node uses a different System ID than the node itself
6. the advertised MTUs match on both sides AND
7. both nodes advertise defined level values AND
8. [
  - i) the node is at level 0 and has no three way adjacencies already to nodes at Highest Adjacency Three-Way level (HAT as defined later in Section 4.2.7.1) with level different than the adjacent node OR
  - ii) the node is not at level 0 and the neighboring node is at level 0 OR
  - iii) both nodes are at level 0 AND both indicate support for Section 4.3.8 OR
  - iv) neither node is at level 0 and the neighboring node is at most one level away].

The rules checking PoD numbering MAY be optionally disregarded by a node if PoD detection is undesirable or has to be ignored. This will not affect the correctness of the protocol except preventing detection of certain miscabling cases.

A node configured with "undefined" PoD membership MUST, after building first northbound three way adjacencies to a node being in a defined PoD, advertise that PoD as part of its LIEs. In case that adjacency is lost, from all available northbound three way adjacencies the node with the highest System ID and defined PoD is chosen. That way the northmost defined PoD value (normally the ToP



nodes) can diffuse southbound towards the leaves "forcing" the PoD value on any node with "undefined" PoD.

LIEs arriving with IPv4 Time to Live (TTL) / IPv6 Hop Limit (HL) larger than 1 MUST be ignored.

A node SHOULD NOT send out LIEs without defined level in the header but in certain scenarios it may be beneficial for trouble-shooting purposes.

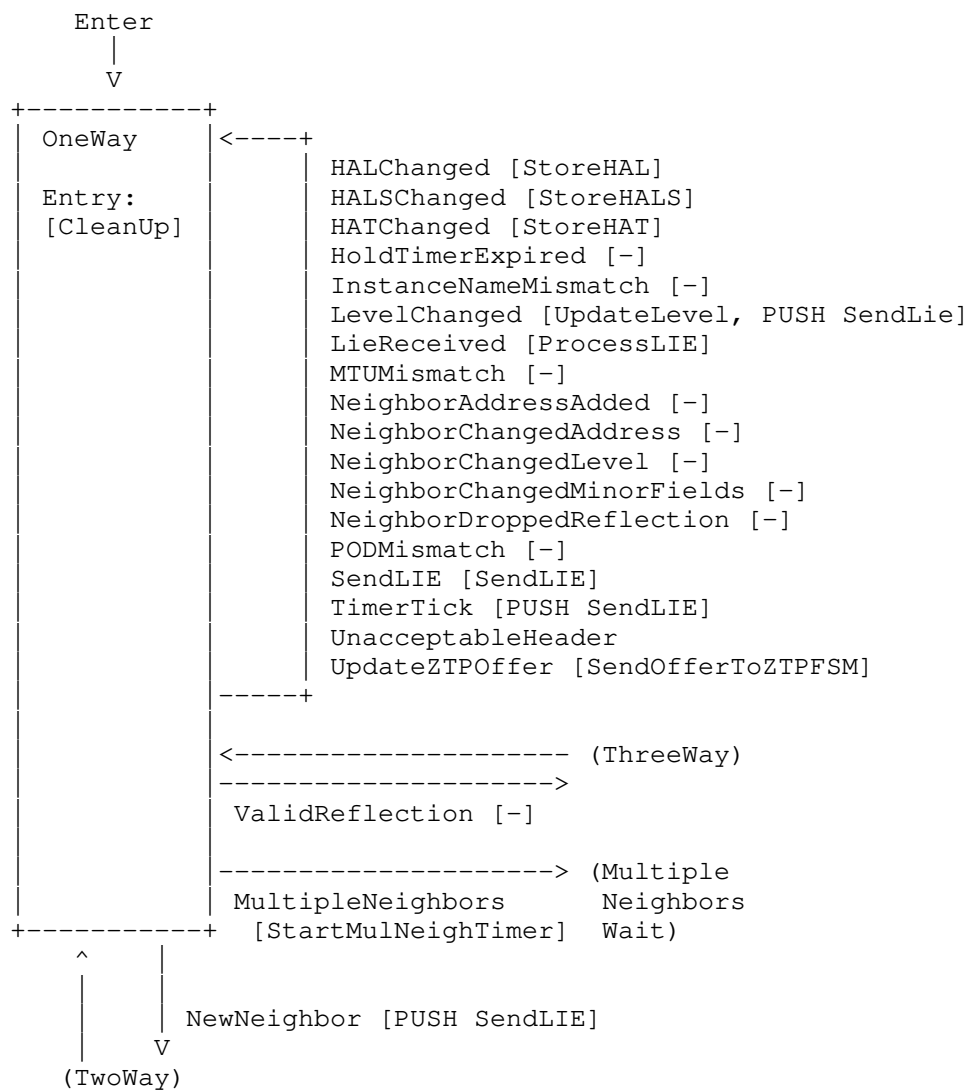
#### 4.2.2.1. LIE FSM

This section specifies the precise, normative LIE FSM and can be omitted unless the reader is pursuing an implementation of the protocol.

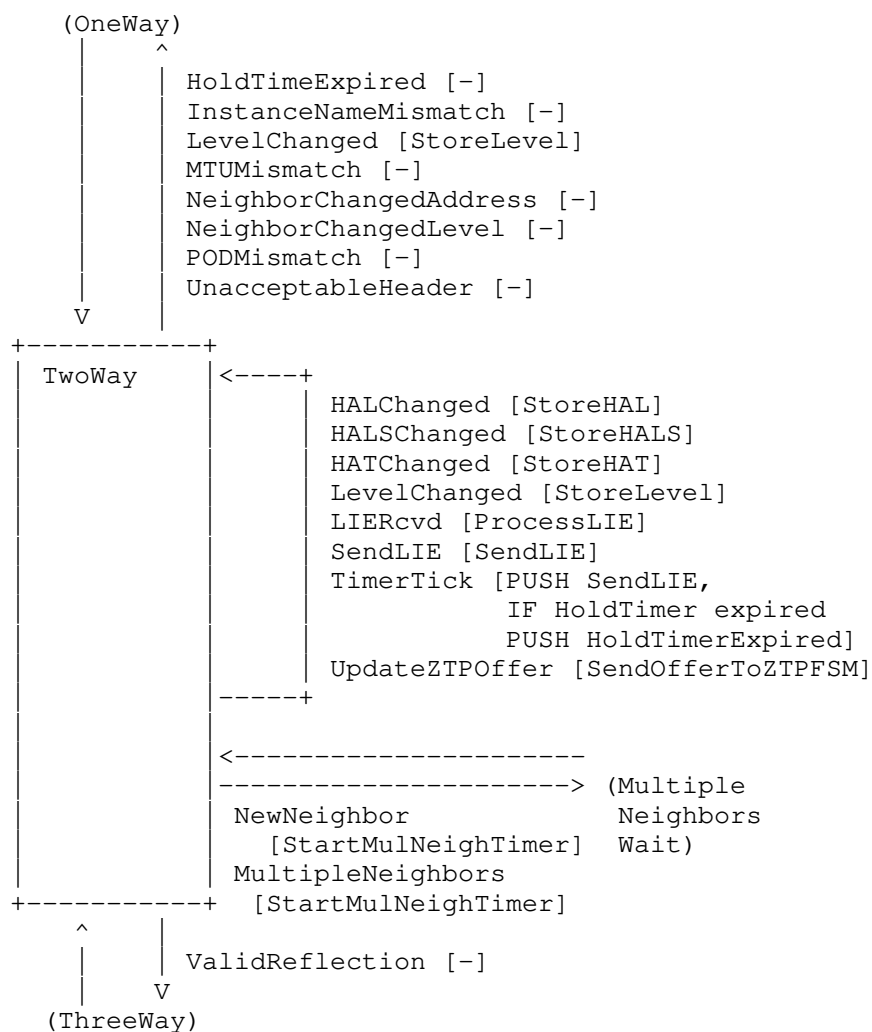
Initial state is 'OneWay'.

Event 'MultipleNeighbors' occurs normally when more than two nodes see each other on the same link or a remote node is quickly reconfigured or rebooted without regressing to 'OneWay' first. Each occurrence of the event SHOULD generate a clear, according notification to help operational deployments.

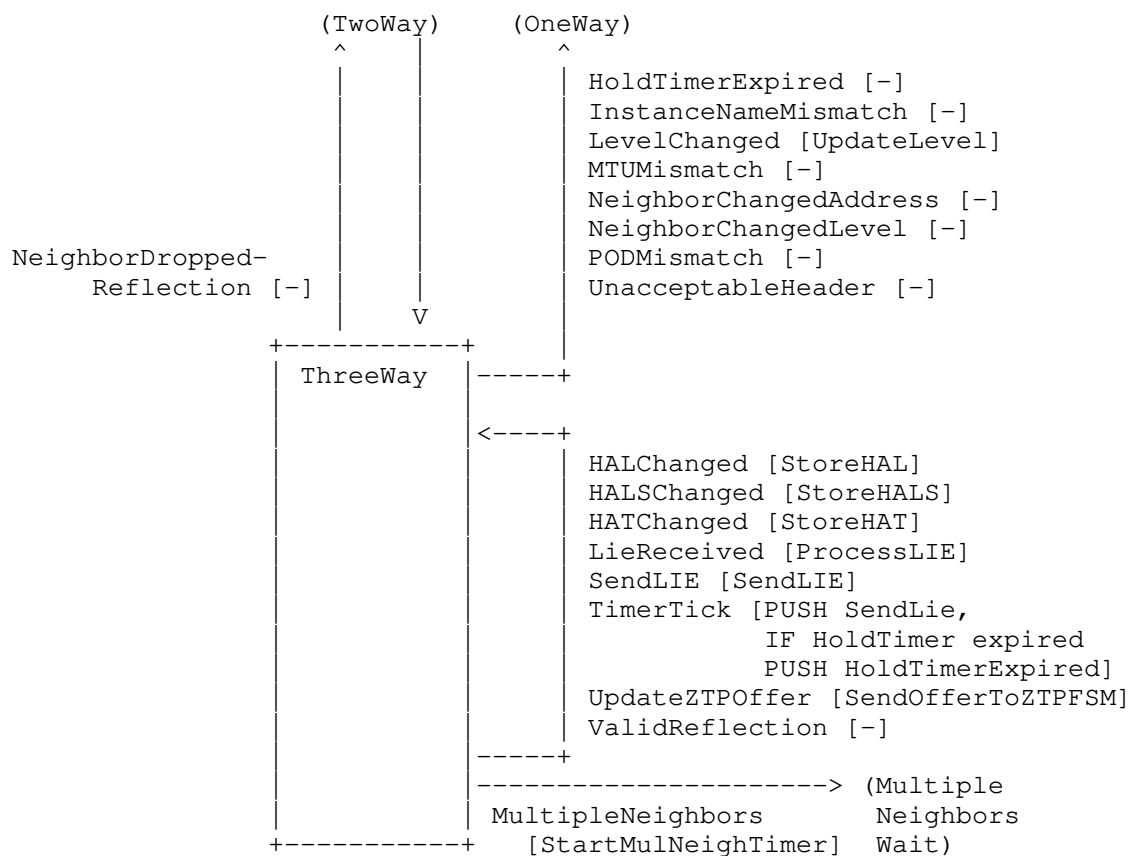
The machine sends LIEs on several transitions to accelerate adjacency bring-up without waiting for the timer tic.



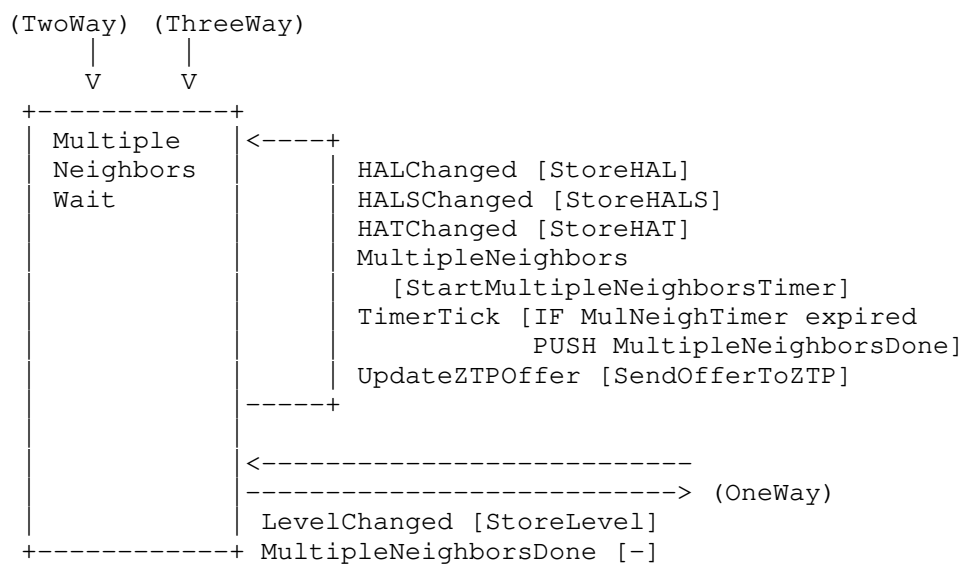
LIE FSM



LIE FSM (continued)



LIE FSM (continued)



LIE FSM (continued)

## Events

- o TimerTick: one second timer tic
- o LevelChanged: node's level has been changed by ZTP or configuration
- o HALChanged: best HAL computed by ZTP has changed
- o HATChanged: HAT computed by ZTP has changed
- o HALSChanged: set of HAL offering systems computed by ZTP has changed
- o LieRcvd: received LIE
- o NewNeighbor: new neighbor parsed
- o ValidReflection: received own reflection from neighbor
- o NeighborDroppedReflection: lost previous own reflection from neighbor
- o NeighborChangedLevel: neighbor changed advertised level
- o NeighborChangedAddress: neighbor changed IP address

- o UnacceptableHeader: unacceptable header seen
- o MTUMismatch: MTU mismatched
- o InstanceNameMismatch: Instance mismatched
- o PODMismatch: Unacceptable PoD seen
- o HoldtimeExpired: adjacency hold down expired
- o MultipleNeighbors: more than one neighbor seen on interface
- o MultipleNeighborsDone: cooldown for multiple neighbors expired
- o SendLie: send a LIE out
- o UpdateZTPOffer: update this node's ZTP offer

#### Actions

on MultipleNeighbors in OneWay finishes in MultipleNeighborsWait:  
start multiple neighbors timer as  $4 * \text{DEFAULT\_LIE\_HOLDTIME}$

on NeighborDroppedReflection in ThreeWay finishes in TwoWay: no  
action

on NeighborDroppedReflection in OneWay finishes in OneWay: no  
action

on PODMismatch in TwoWay finishes in OneWay: no action

on NewNeighbor in TwoWay finishes in MultipleNeighborsWait: PUSH  
SendLie event

on LieRcvd in OneWay finishes in OneWay: PROCESS\_LIE

on UnacceptableHeader in ThreeWay finishes in OneWay: no action

on UpdateZTPOffer in TwoWay finishes in TwoWay: send offer to ZTP  
FSM

on NeighborChangedAddress in ThreeWay finishes in OneWay: no  
action

on HALChanged in MultipleNeighborsWait finishes in  
MultipleNeighborsWait: store new HAL

on NeighborChangedAddress in TwoWay finishes in OneWay: no action

on MultipleNeighbors in TwoWay finishes in MultipleNeighborsWait: start multiple neighbors timer as 4 \* DEFAULT\_LIE\_HOLDTIME

on LevelChanged in ThreeWay finishes in OneWay: update level with event value

on LieRcvd in ThreeWay finishes in ThreeWay: PROCESS\_LIE

on ValidReflection in OneWay finishes in ThreeWay: no action

on NeighborChangedLevel in TwoWay finishes in OneWay: no action

on MultipleNeighbors in ThreeWay finishes in MultipleNeighborsWait: start multiple neighbors timer as 4 \* DEFAULT\_LIE\_HOLDTIME

on InstanceNameMismatch in OneWay finishes in OneWay: no action

on NewNeighbor in OneWay finishes in TwoWay: PUSH SendLie event

on UpdateZTPOffer in OneWay finishes in OneWay: send offer to ZTP FSM

on UpdateZTPOffer in ThreeWay finishes in ThreeWay: send offer to ZTP FSM

on MTUMismatch in ThreeWay finishes in OneWay: no action

on TimerTick in OneWay finishes in OneWay: PUSH SendLie event

on SendLie in TwoWay finishes in TwoWay: SEND\_LIE

on ValidReflection in ThreeWay finishes in ThreeWay: no action

on InstanceNameMismatch in TwoWay finishes in OneWay: no action

on HoldtimeExpired in OneWay finishes in OneWay: no action

on TimerTick in ThreeWay finishes in ThreeWay: PUSH SendLie event, if holdtime expired PUSH HoldtimeExpired event

on HALChanged in TwoWay finishes in TwoWay: store new HAL

on HoldtimeExpired in ThreeWay finishes in OneWay: no action

on HALSChanged in TwoWay finishes in TwoWay: store HALS

on HALSChanged in ThreeWay finishes in ThreeWay: store HALS

on ValidReflection in TwoWay finishes in ThreeWay: no action

on MultipleNeighborsDone in MultipleNeighborsWait finishes in OneWay: no action

on NeighborAddressAdded in OneWay finishes in OneWay: no action

on TimerTick in MultipleNeighborsWait finishes in MultipleNeighborsWait: decrement MultipleNeighbors timer, if expired PUSH MultipleNeighborsDone

on MTUMismatch in OneWay finishes in OneWay: no action

on MultipleNeighbors in MultipleNeighborsWait finishes in MultipleNeighborsWait: start multiple neighbors timer as 4 \* DEFAULT\_LIE\_HOLDTIME

on LieRcvd in TwoWay finishes in TwoWay: PROCESS\_LIE

on HATChanged in MultipleNeighborsWait finishes in MultipleNeighborsWait: store HAT

on HoldtimeExpired in TwoWay finishes in OneWay: no action

on NeighborChangedLevel in ThreeWay finishes in OneWay: no action

on LevelChanged in OneWay finishes in OneWay: update level with event value, PUSH SendLie event

on SendLie in OneWay finishes in OneWay: SEND\_LIE

on HATChanged in OneWay finishes in OneWay: store HAT

on LevelChanged in TwoWay finishes in TwoWay: update level with event value

on HATChanged in TwoWay finishes in TwoWay: store HAT

on PODMismatch in ThreeWay finishes in OneWay: no action

on LevelChanged in MultipleNeighborsWait finishes in OneWay: update level with event value

on UnacceptableHeader in TwoWay finishes in OneWay: no action

on NeighborChangedLevel in OneWay finishes in OneWay: no action

on InstanceNameMismatch in ThreeWay finishes in OneWay: no action



on HATChanged in ThreeWay finishes in ThreeWay: store HAT

on HALChanged in OneWay finishes in OneWay: store new HAL

on UnacceptableHeader in OneWay finishes in OneWay: no action

on HALChanged in ThreeWay finishes in ThreeWay: store new HAL

on UpdateZTPOffer in MultipleNeighborsWait finishes in MultipleNeighborsWait: send offer to ZTP FSM

on NeighborChangedMinorFields in OneWay finishes in OneWay: no action

on NeighborChangedAddress in OneWay finishes in OneWay: no action

on MTUMismatch in TwoWay finishes in OneWay: no action

on PODMismatch in OneWay finishes in OneWay: no action

on SendLie in ThreeWay finishes in ThreeWay: SEND\_LIE

on TimerTick in TwoWay finishes in TwoWay: PUSH SendLie event, if holdtime expired PUSH HoldtimeExpired event

on HALSChanged in OneWay finishes in OneWay: store HALS

on HALSChanged in MultipleNeighborsWait finishes in MultipleNeighborsWait: store HALS

on Entry into OneWay: CLEANUP

Following words are used for well known procedures:

1. PUSH Event: pushes an event to be executed by the FSM upon exit of this action
2. CLEANUP: neighbor MUST be reset to unknown
3. SEND\_LIE: create a new LIE packet
  1. reflecting the neighbor if known and valid and
  2. setting the necessary 'not\_a\_ztp\_offer' variable if level was derived from last known neighbor on this interface and
  3. setting 'you\_are\_flood\_repeater' to computed value

4. PROCESS\_LIE:
  1. if lie has wrong major version OR our own system ID or invalid system ID then CLEANUP else
  2. if lie has non matching MTUs then CLEANUP, PUSH UpdateZTPOffer, PUSH MTUMismatch else
  3. if PoD rules do not allow adjacency forming then CLEANUP, PUSH PODMismatch, PUSH MTUMismatch else
  4. if lie has undefined level OR my level is undefined OR this node is leaf and remote level lower than HAT OR (lie's level is not leaf AND its difference is more than one from my level) then CLEANUP, PUSH UpdateZTPOffer, PUSH UnacceptableHeader else
  5. PUSH UpdateZTPOffer, construct temporary new neighbor structure with values from lie, if no current neighbor exists then set neighbor to new neighbor, PUSH NewNeighbor event, CHECK\_THREE\_WAY else
    1. if current neighbor system ID differs from lie's system ID then PUSH MultipleNeighbors else
    2. if current neighbor stored level differs from lie's level then PUSH NeighborChangedLevel else
    3. if current neighbor stored IPv4/v6 address differs from lie's address then PUSH NeighborChangedAddress else
    4. if any of neighbor's flood address port, name, local linkid changed then PUSH NeighborChangedMinorFields and
    5. CHECK\_THREE\_WAY
5. CHECK\_THREE\_WAY: if current state is one-way do nothing else
  1. if lie packet does not contain neighbor then if current state is three-way then PUSH NeighborDroppedReflection else
  2. if packet reflects this system's ID and local port and state is three-way then PUSH event ValidReflection else PUSH event MultipleNeighbors

### 4.2.3. Topology Exchange (TIE Exchange)

#### 4.2.3.1. Topology Information Elements

Topology and reachability information in RIFT is conveyed by the means of TIEs which have good amount of commonalities with LSAs in OSPF.

The TIE exchange mechanism uses the port indicated by each node in the LIE exchange and the interface on which the adjacency has been formed as destination. It SHOULD use TTL of 1 as well and set inter-network control precedence on according packets.

TIEs contain sequence numbers, lifetimes and a type. Each type has ample identifying number space and information is spread across possibly many TIEs of a certain type by the means of a hash function that a node or deployment can individually determine. One extreme design choice is a prefix per TIE which leads to more BGP-like behavior where small increments are only advertised on route changes vs. deploying with dense prefix packing into few TIEs leading to more traditional IGP trade-off with fewer TIEs. An implementation may even rehash prefix to TIE mapping at any time at the cost of significant amount of re-advertisements of TIEs.

More information about the TIE structure can be found in the schema in Appendix B.

#### 4.2.3.2. South- and Northbound Representation

A central concept of RIFT is that each node represents itself differently depending on the direction in which it is advertising information. More precisely, a spine node represents two different databases over its adjacencies depending whether it advertises TIEs to the north or to the south/sideways. We call those differing TIE databases either south- or northbound (South TIEs and North TIEs) depending on the direction of distribution.

The North TIEs hold all of the node's adjacencies and local prefixes while the South TIEs hold only all of the node's adjacencies, the default prefix with necessary disaggregated prefixes and local prefixes. We will explain this in detail further in Section 4.2.5.

The TIE types are mostly symmetric in both directions and Table 2 provides a quick reference to main TIE types including direction and their function.

TIE-Type	Content
Node North TIE	node properties and adjacencies
Node South TIE	same content as node North TIE
Prefix North TIE	contains nodes' directly reachable prefixes
Prefix South TIE	contains originated defaults and directly reachable prefixes
Positive Disaggregation South TIE	contains disaggregated prefixes
Negative Disaggregation South TIE	contains special, negatively disaggregated prefixes to support multi-plane designs
External Prefix North TIE	contains external prefixes
Key-Value North TIE	contains nodes northbound KVs
Key-Value South TIE	contains nodes southbound KVs

Table 2: TIE Types

As an example illustrating a databases holding both representations, consider the topology in Figure 2 with the optional link between spine 111 and spine 112 (so that the flooding on an East-West link can be shown). This example assumes unnumbered interfaces. First, here are the TIEs generated by some nodes. For simplicity, the key value elements which may be included in their South TIEs or North TIEs are not shown.

ToF 21 South TIEs:

Node South TIE:

```
NodeElement(level=2, neighbors((Spine 111, level 1, cost 1),
(Spine 112, level 1, cost 1), (Spine 121, level 1, cost 1),
(Spine 122, level 1, cost 1)))
```

Prefix South TIE:

```
SouthPrefixesElement(prefixes(0/0, cost 1), (::/0, cost 1))
```

Spine 111 South TIEs:

Node South TIE:

```
NodeElement(level=1, neighbors((ToF 21, level 2, cost 1,
    links(...)),
    (ToF 22, level 2, cost 1, links(...)),
    (Spine 112, level 1, cost 1, links(...)),
    (Leaf111, level 0, cost 1, links(...)),
    (Leaf112, level 0, cost 1, links(...))))
```

Prefix South TIE:

```
SouthPrefixesElement(prefixes(0/0, cost 1), (::/0, cost 1))
```

Spine 111 North TIEs:

Node North TIE:

```
NodeElement(level=1,
    neighbors((ToF 21, level 2, cost 1, links(...)),
    (ToF 22, level 2, cost 1, links(...)),
    (Spine 112, level 1, cost 1, links(...)),
    (Leaf111, level 0, cost 1, links(...)),
    (Leaf112, level 0, cost 1, links(...))))
```

Prefix North TIE:

```
NorthPrefixesElement(prefixes(Spine 111.loopback)
```

Spine 121 South TIEs:

Node South TIE:

```
NodeElement(level=1, neighbors((ToF 21, level 2, cost 1),
    (ToF 22, level 2, cost 1), (Leaf121, level 0, cost 1),
    (Leaf122, level 0, cost 1)))
```

Prefix South TIE:

```
SouthPrefixesElement(prefixes(0/0, cost 1), (::/0, cost 1))
```

Spine 121 North TIEs:

Node North TIE:

```
NodeElement(level=1,
    neighbors((ToF 21, level 2, cost 1, links(...)),
    (ToF 22, level 2, cost 1, links(...)),
    (Leaf121, level 0, cost 1, links(...)),
    (Leaf122, level 0, cost 1, links(...))))
```

Prefix North TIE:

```
NorthPrefixesElement(prefixes(Spine 121.loopback)
```

Leaf112 North TIEs:

Node North TIE:

```
NodeElement(level=0,
    neighbors((Spine 111, level 1, cost 1, links(...)),
    (Spine 112, level 1, cost 1, links(...))))
```

Prefix North TIE:

```
NorthPrefixesElement(prefixes(Leaf112.loopback, Prefix112,
    Prefix_MH))
```

Figure 14: Example TIES Generated in a 2 Level Spine-and-Leaf Topology

It may be here not necessarily obvious why the node South TIES contain all the adjacencies of the according node. This will be necessary for algorithms given in Section 4.2.3.9 and Section 4.3.6.

#### 4.2.3.3. Flooding

The mechanism used to distribute TIES is the well-known (albeit modified in several respects to take advantage of fat tree topology) flooding mechanism used by today's link-state protocols. Although flooding is initially more demanding to implement it avoids many problems with update style used in diffused computation such as distance vector protocols. Since flooding tends to present an unscalable burden in large, densely meshed topologies (fat trees being unfortunately such a topology) we provide as solution close to optimal global flood reduction and load balancing optimization in Section 4.2.3.9.

As described before, TIES themselves are transported over UDP with the ports indicated in the LIE exchanges and using the destination address on which the LIE adjacency has been formed. For unnumbered IPv4 interfaces same considerations apply as in equivalent OSPF case.

##### 4.2.3.3.1. Normative Flooding Procedures

On reception of a TIE with an undefined level value in the packet header the node SHOULD issue a warning and indiscriminately discard the packet.

This section specifies the precise, normative flooding mechanism and can be omitted unless the reader is pursuing an implementation of the protocol.

Flooding Procedures are described in terms of a flooding state of an adjacency and resulting operations on it driven by packet arrivals. The FSM itself has basically just a single state and is not well suited to represent the behavior. An implementation MUST behave on the wire in the same way as the provided normative procedures of this paragraph.

RIFT does not specify any kind of flood rate limiting since such specifications always assume particular points in available technology speeds and feeds and those points are shifting at faster and faster rate (speed of light holding for the moment). The encoded packets provide hints to react accordingly to losses or overruns.

Flooding of all according topology exchange elements SHOULD be performed at highest feasible rate whereas the rate of transmission MUST be throttled by reacting to adequate features of the system such as e.g. queue lengths or congestion indications in the protocol packets.

A node SHOULD NOT send out any topology information elements if the adjacency is not in a "three-way" state. No further tightening of this rule is possible due to possible link buffering and re-ordering of LIEs and TIEs/TIDEs/TIREs.

A node MUST drop any received TIEs/TIDEs/TIREs unless it is in three-way state.

TIDEs and TIREs MUST NOT be re-flooded the way TIEs of other nodes MUST be always generated by the node itself and cross only to the neighboring node.

#### 4.2.3.3.1.1. FloodState Structure per Adjacency

The structure contains conceptually the following elements. The word collection or queue indicates a set of elements that can be iterated:

TIES\_TX: Collection containing all the TIEs to transmit on the adjacency.

TIES\_ACK: Collection containing all the TIEs that have to be acknowledged on the adjacency.

TIES\_REQ: Collection containing all the TIE headers that have to be requested on the adjacency.

TIES\_RTX: Collection containing all TIEs that need retransmission with the according time to retransmit.

Following words are used for well known procedures operating on this structure:

TIE Describes either a full RIFT TIE or accordingly just the 'TIEHeader' or 'TIEID'. The according meaning is unambiguously contained in the context of the algorithm.

is\_flood\_reduced(TIE): returns whether a TIE can be flood reduced or not.

is\_tide\_entry\_filtered(TIE): returns whether a header should be propagated in TIDE according to flooding scopes.

`is_request_filtered(TIE)`: returns whether a TIE request should be propagated to neighbor or not according to flooding scopes.

`is_flood_filtered(TIE)`: returns whether a TIE requested be flooded to neighbor or not according to flooding scopes.

`try_to_transmit_tie(TIE)`:

A. if not `is_flood_filtered(TIE)` then

1. remove TIE from TIES\_RTX if present

2. if TIE" with same key on TIES\_ACK then

a. if TIE" same or newer than TIE do nothing else

b. remove TIE" from TIES\_ACK and add TIE to TIES\_TX

3. else insert TIE into TIES\_TX

`ack_tie(TIE)`: remove TIE from all collections and then insert TIE into TIES\_ACK.

`tie_been_acked(TIE)`: remove TIE from all collections.

`remove_from_all_queues(TIE)`: same as `'tie_been_acked'`.

`request_tie(TIE)`: if not `is_request_filtered(TIE)` then `remove_from_all_queues(TIE)` and add to TIES\_REQ.

`move_to_rtx_list(TIE)`: remove TIE from TIES\_TX and then add to TIES\_RTX using TIE retransmission interval.

`clear_requests(TIEs)`: remove all TIEs from TIES\_REQ.

`bump_own_tie(TIE)`: for self-originated TIE originate an empty or re-generate with version number higher then the one in TIE.

The collection SHOULD be served with following priorities if the system cannot process all the collections in real time:

Elements on TIES\_ACK should be processed with highest priority

TIES\_TX

TIES\_REQ and TIES\_RTX



## 4.2.3.3.1.2. TIDEs

`TIEID` and `TIEHeader` space forms a strict total order (modulo incomparable sequence numbers in the very unlikely event that can occur if a TIE is "stuck" in a part of a network while the originator reboots and reissues TIEs many times to the point its sequence# rolls over and forms incomparable distance to the "stuck" copy) which implies that a comparison relation is possible between two elements. With that it is implicitly possible to compare TIEs, TIEHeaders and TIEIDs to each other whereas the shortest viable key is always implied.

When generating and sending TIDEs an implementation SHOULD ensure that enough bandwidth is left to send elements of Floodstate structure.

## 4.2.3.3.1.2.1. TIDE Generation

As given by timer constant, periodically generate TIDEs by:

NEXT\_TIDE\_ID: ID of next TIE to be sent in TIDE.

TIDE\_START: Begin of TIDE packet range.

- a. NEXT\_TIDE\_ID = MIN\_TIEID
- b. while NEXT\_TIDE\_ID not equal to MAX\_TIEID do
  1. TIDE\_START = NEXT\_TIDE\_ID
  2. HEADERS = At most TIRDEs\_PER\_PKT headers in TIEDB starting at NEXT\_TIDE\_ID or higher that SHOULD be filtered by is\_tide\_entry\_filtered and MUST either have a lifetime left > 0 or have no content
  3. if HEADERS is empty then START = MIN\_TIEID else START = first element in HEADERS
  4. if HEADERS' size less than TIRDEs\_PER\_PKT then END = MAX\_TIEID else END = last element in HEADERS
  5. send sorted HEADERS as TIDE setting START and END as its range
  6. NEXT\_TIDE\_ID = END

The constant `TIRDES\_PER\_PKT` SHOULD be generated and used by the implementation to limit the amount of TIE headers per TIDE so the sent TIDE PDU does not exceed interface MTU.

TIDE PDUs SHOULD be spaced on sending to prevent packet drops.

#### 4.2.3.3.1.2.2. TIDE Processing

On reception of TIDEs the following processing is performed:

TXKEYS: Collection of TIE Headers to be send after processing of the packet

REQKEYS: Collection of TIEIDs to be requested after processing of the packet

CLEARKEYS: Collection of TIEIDs to be removed from flood state queues

LASTPROCESSED: Last processed TIEID in TIDE

DBTIE: TIE in the LSDB if found

- a. LASTPROCESSED = TIDE.start\_range
- b. for every HEADER in TIDE do
  1. DBTIE = find HEADER in current LSDB
  2. if HEADER < LASTPROCESSED then report error and reset adjacency and return
  3. put all TIEs in LSDB where (TIE.HEADER > LASTPROCESSED and TIE.HEADER < HEADER) into TXKEYS
  4. LASTPROCESSED = HEADER
  5. if DBTIE not found then
    - I) if originator is this node then bump\_own\_tie
    - II) else put HEADER into REQKEYS
  6. if DBTIE.HEADER < HEADER then
    - I) if originator is this node then bump\_own\_tie else

- i. if this is a North TIE header from a northbound neighbor then override DBTIE in LSDB with HEADER
  - ii. else put HEADER into REQKEYS
- 7. if DBTIE.HEADER > HEADER then put DBTIE.HEADER into TXKEYS
- 8. if DBTIE.HEADER = HEADER then
  - I) if DBTIE has content already then put DBTIE.HEADER into CLEARKEYS
  - II) else put HEADER into REQKEYS
- c. put all TIEs in LSDB where (TIE.HEADER > LASTPROCESSED and TIE.HEADER <= TIE.end\_range) into TXKEYS
- d. for all TIEs in TXKEYS try\_to\_transmit\_tie(TIE)
- e. for all TIEs in REQKEYS request\_tie(TIE)
- f. for all TIEs in CLEARKEYS remove\_from\_all\_queues(TIE)

#### 4.2.3.3.1.3. TIREs

##### 4.2.3.3.1.3.1. TIRE Generation

Elements from both TIES\_REQ and TIES\_ACK MUST be collected and sent out as fast as feasible as TIREs. When sending TIREs with elements from TIES\_REQ the 'lifetime' field MUST be set to 0 to force reflooding from the neighbor even if the TIEs seem to be same.

##### 4.2.3.3.1.3.2. TIRE Processing

On reception of TIREs the following processing is performed:

TXKEYS: Collection of TIE Headers to be send after processing of the packet

REQKEYS: Collection of TIEIDs to be requested after processing of the packet

ACKKEYS: Collection of TIEIDs that have been acked

DBTIE: TIE in the LSDB if found

- a. for every HEADER in TIRE do

1. DBTIE = find HEADER in current LSDB
  2. if DBTIE not found then do nothing
  3. if DBTIE.HEADER < HEADER then put HEADER into REQKEYS
  4. if DBTIE.HEADER > HEADER then put DBTIE.HEADER into TXKEYS
  5. if DBTIE.HEADER = HEADER then put DBTIE.HEADER into ACKKEYS
  - b. for all TIES in TXKEYS try\_to\_transmit\_tie(TIE)
  - c. for all TIES in REQKEYS request\_tie(TIE)
  - d. for all TIES in ACKKEYS tie\_been\_acked(TIE)
- 4.2.3.3.1.4. TIES Processing on Flood State Adjacency

On reception of TIES the following processing is performed:

ACKTIE: TIE to acknowledge

TXTIE: TIE to transmit

DBTIE: TIE in the LSDB if found

- a. DBTIE = find TIE in current LSDB
- b. if DBTIE not found then
  1. if originator is this node then bump\_own\_tie with a short remaining lifetime
  2. else insert TIE into LSDB and ACKTIE = TIE
- else
  1. if DBTIE.HEADER = TIE.HEADER then
    - i. if DBTIE has content already then ACKTIE = TIE
    - ii. else process like the "DBTIE.HEADER < TIE.HEADER" case
  2. if DBTIE.HEADER < TIE.HEADER then
    - i. if originator is this node then bump\_own\_tie
    - ii. else insert TIE into LSDB and ACKTIE = TIE

3. if DBTIE.HEADER > TIE.HEADER then
  - i. if DBTIE has content already then TXTIE = DBTIE
  - ii. else ACKTIE = DBTIE
- c. if TXTIE is set then try\_to\_transmit\_tie(TXTIE)
- d. if ACKTIE is set then ack\_tie(TIE)

#### 4.2.3.3.1.5. TIEs Processing When LSDB Received Newer Version on Other Adjacencies

The Link State Database can be considered to be a switchboard that does not need any flooding procedures but can be given new versions of TIEs by a peer. Consecutively, a peer receives from the LSDB newer versions of TIEs received by other peers and processes them (without any filtering) just like receiving TIEs from its remote peer. This publisher model can be implemented in many ways.

#### 4.2.3.3.1.6. Sending TIEs

On a periodic basis all TIEs with lifetime left > 0 MUST be sent out on the adjacency, removed from TIES\_TX list and requeued onto TIES\_RTX list.

#### 4.2.3.4. TIE Flooding Scopes

In a somewhat analogous fashion to link-local, area and domain flooding scopes, RIFT defines several complex "flooding scopes" depending on the direction and type of TIE propagated.

Every North TIE is flooded northbound, providing a node at a given level with the complete topology of the Clos or Fat Tree network that is reachable southwards of it, including all specific prefixes. This means that a packet received from a node at the same or lower level whose destination is covered by one of those specific prefixes will be routed directly towards the node advertising that prefix rather than sending the packet to a node at a higher level.

A node's Node South TIEs, consisting of all node's adjacencies and prefix South TIEs limited to those related to default IP prefix and disaggregated prefixes, are flooded southbound in order to allow the nodes one level down to see connectivity of the higher level as well as reachability to the rest of the fabric. In order to allow an E-W disconnected node in a given level to receive the South TIEs of other nodes at its level, every \*NODE\* South TIE is "reflected" northbound to level from which it was received. It should be noted that East-

West links are included in South TIE flooding (except at ToF level); those TIEs need to be flooded to satisfy algorithms in Section 4.2.4. In that way nodes at same level can learn about each other without a lower level, e.g. in case of leaf level. The precise, normative flooding scopes are given in Table 3. Those rules govern as well what SHOULD be included in TIEs on the adjacency. Again, East-West flooding scopes are identical to South flooding scopes except in case of ToF East-West links (rings) which are basically performing northbound flooding.

Node South TIE "south reflection" allows to support positive disaggregation on failures describes in Section 4.2.5 and flooding reduction in Section 4.2.3.9.

Type / Direction	South	North	East-West
node South TIE	flood if level of originator is equal to this node	flood if level of originator is higher than this node	flood only if this node is not ToF
non-node South TIE	flood self-originated only	flood only if neighbor is originator of TIE	flood only if self-originated and this node is not ToF
all North TIEs	never flood	flood always	flood only if this node is ToF
TIDE	include at least all non-self originated North TIE headers and self-originated South TIE headers and node South TIEs of nodes at same level	include at least all node South TIEs and all South TIEs originated by peer and all North TIEs	if this node is ToF then include all North TIEs, otherwise only self-originated TIEs
TIRE as Request	request all North TIEs and all peer's self-originated TIEs and all node South TIEs	request all South TIEs	if this node is ToF then apply North scope rules, otherwise South scope rules
TIRE as Ack	Ack all received TIEs	Ack all received TIEs	Ack all received TIEs

Table 3: Normative Flooding Scopes

If the TIDE includes additional TIE headers beside the ones specified, the receiving neighbor must apply according filter to the received TIDE strictly and MUST NOT request the extra TIE headers that were not allowed by the flooding scope rules in its direction.

As an example to illustrate these rules, consider using the topology in Figure 2, with the optional link between spine 111 and spine 112, and the associated TIEs given in Figure 14. The flooding from particular nodes of the TIEs is given in Table 4.

Router floods to	Neighbor	TIEs
Leaf111	Spine 112	Leaf111 North TIEs, Spine 111 node South TIE
Leaf111	Spine 111	Leaf111 North TIEs, Spine 112 node South TIE
Spine 111	Leaf111	Spine 111 South TIEs
Spine 111	Leaf112	Spine 111 South TIEs
Spine 111	Spine 112	Spine 111 South TIEs
Spine 111	ToF 21	Spine 111 North TIEs, Leaf111 North TIEs, Leaf112 North TIEs, ToF 22 node South TIE
Spine 111	ToF 22	Spine 111 North TIEs, Leaf111 North TIEs, Leaf112 North TIEs, ToF 21 node South TIE
...	...	...
ToF 21	Spine 111	ToF 21 South TIEs
ToF 21	Spine 112	ToF 21 South TIEs
ToF 21	Spine 121	ToF 21 South TIEs
ToF 21	Spine 122	ToF 21 South TIEs
...	...	...

Table 4: Flooding some TIEs from example topology

#### 4.2.3.5. 'Flood Only Node TIEs' Bit

RIFT includes an optional ECN (Explicit Congestion Notification) mechanism to prevent "flooding inrush" on restart or bring-up with many southbound neighbors. A node MAY set on its LIEs the according bit to indicate to the neighbor that it should temporarily flood node TIEs only to it. It SHOULD only set it in the southbound direction. The receiving node SHOULD accommodate the request to lessen the



flooding load on the affected node if south of the sender and SHOULD ignore the bit if northbound.

Obviously this mechanism is most useful in southbound direction. The distribution of node TIEs guarantees correct behavior of algorithms like disaggregation or default route origination. Furthermore though, the use of this bit presents an inherent trade-off between processing load and convergence speed since suppressing flooding of northbound prefixes from neighbors will lead to blackholes.

#### 4.2.3.6. Initial and Periodic Database Synchronization

The initial exchange of RIFT is modeled after ISIS with TIDE being equivalent to CSNP and TIRE playing the role of PSNP. The content of TIDEs and TIREs is governed by Table 3.

#### 4.2.3.7. Purging and Roll-Overs

When a node exits the network, if "unpurged", residual stale TIEs may exist in the network until their lifetimes expire (which in case of RIFT is by default a rather long period to prevent ongoing re-origination of TIEs in very large topologies). RIFT does however not have a "purging mechanism" in the traditional sense based on sending specialized "purge" packets. In other routing protocols such mechanism has proven to be complex and fragile based on many years of experience. RIFT simply issues a new, empty version of the TIE with a short lifetime and relies on each node to age out and delete such TIE copy independently. Abundant amounts of memory are available today even on low-end platforms and hence keeping those relatively short-lived extra copies for a while is acceptable. The information will age out and in the meantime all computations will deliver correct results if a node leaves the network due to the new information distributed by its adjacent nodes breaking bi-directional connectivity checks in different computations.

Once a RIFT node issues a TIE with an ID, it SHOULD preserve the ID as long as feasible (also when the protocol restarts), even if the TIE loses all content. The re-advertisement of empty TIE fulfills the purpose of purging any information advertised in previous versions. The originator is free to not re-originate the according empty TIE again or originate an empty TIE with relatively short lifetime to prevent large number of long-lived empty stubs polluting the network. Each node MUST timeout and clean up the according empty TIEs independently.

Upon restart a node MUST, as any link-state implementation, be prepared to receive TIEs with its own system ID and supersede them with equivalent, newly generated, empty TIEs with a higher sequence

number. As above, the lifetime can be relatively short since it only needs to exceed the necessary propagation and processing delay by all the nodes that are within the TIE's flooding scope.

TIE sequence numbers are rolled over using the method described in Appendix A. First sequence number of any spontaneously originated TIE (i.e. not originated to override a detected older copy in the network) MUST be a reasonably unpredictable random number in the interval  $[0, 2^{30}-1]$  which will prevent otherwise identical TIE headers to remain "stuck" in the network with content different from TIE originated after reboot. In traditional link-state protocols this is delegated to a 16-bit checksum on packet content. RIFT avoids this design due to the CPU burden presented by computation of such checksums and additional complications tied to the fact that the checksum must be "patched" into the packet after the computation, a difficult proposition in binary hand-crafted formats already and highly incompatible with model-based, serialized formats. The sequence number space is hence consciously chosen to be 64-bits wide to make the occurrence of a TIE with same sequence number but different content as much or even more unlikely than the checksum method. To emulate the "checksum behavior" an implementation could e.g. choose to compute 64-bit checksum over the packet content and use that as first sequence number after reboot.

#### 4.2.3.8. Southbound Default Route Origination

Under certain conditions nodes issue a default route in their South Prefix TIEs with costs as computed in Section 4.3.6.1.

A node X that

1. is NOT overloaded AND
2. has southbound or East-West adjacencies

originates in its south prefix TIE such a default route IIF

1. all other nodes at X's' level are overloaded OR
2. all other nodes at X's' level have NO northbound adjacencies OR
3. X has computed reachability to a default route during N-SPF.

The term "all other nodes at X's' level" describes obviously just the nodes at the same level in the PoD with a viable lower level (otherwise the node South TIEs cannot be reflected and the nodes in e.g. PoD 1 and PoD 2 are "invisible" to each other).

A node originating a southbound default route MUST install a default discard route if it did not compute a default route during N-SPF.

#### 4.2.3.9. Northbound TIE Flooding Reduction

Section 1.4 of the Optimized Link State Routing Protocol [RFC3626] (OLSR) introduces the concept of a "multipoint relay" (MPR) that minimize the overhead of flooding messages in the network by reducing redundant retransmissions in the same region.

A similar technique is applied to RIFT to control northbound flooding. Important observations first:

1. a node MUST flood self-originated North TIEs to all the reachable nodes at the level above which we call the node's "parents";
2. it is typically not necessary that all parents reflowd the North TIEs to achieve a complete flooding of all the reachable nodes two levels above which we choose to call the node's "grandparents";
3. to control the volume of its flooding two hops North and yet keep it robust enough, it is advantageous for a node to select a subset of its parents as "Flood Repeaters" (FRs), which combined together deliver two or more copies of its flooding to all of its parents, i.e. the originating node's grandparents;
4. nodes at the same level do NOT have to agree on a specific algorithm to select the FRs, but overall load balancing should be achieved so that different nodes at the same level should tend to select different parents as FRs;
5. there are usually many solutions to the problem of finding a set of FRs for a given node; the problem of finding the minimal set is (similar to) a NP-Complete problem and a globally optimal set may not be the minimal one if load-balancing with other nodes is an important consideration;
6. it is expected that there will be often sets of equivalent nodes at a level L, defined as having a common set of parents at L+1. Applying this observation at both L and L+1, an algorithm may attempt to split the larger problem in a sum of smaller separate problems;
7. it is another expectation that there will be from time to time a broken link between a parent and a grandparent, and in that case the parent is probably a poor FR due to its lower reliability. An algorithm may attempt to eliminate parents with broken

northbound adjacencies first in order to reduce the number of FRs. Albeit it could be argued that relying on higher fanout FRs will slow flooding due to higher replication load reliability of FR's links seems to be a more pressing concern.

In a fully connected Clos Network, this means that a node selects one arbitrary parent as FR and then a second one for redundancy. The computation can be kept relatively simple and completely distributed without any need for synchronization amongst nodes. In a "PoD" structure, where the Level L+2 is partitioned in silos of equivalent grandparents that are only reachable from respective parents, this means treating each silo as a fully connected Clos Network and solve the problem within the silo.

In terms of signaling, a node has enough information to select its set of FRs; this information is derived from the node's parents' Node South TIEs, which indicate the parent's reachable northbound adjacencies to its own parents, i.e. the node's grandparents. A node may send a LIE to a northbound neighbor with the optional boolean field 'you\_are\_flood\_repeater' set to false, to indicate that the northbound neighbor is not a flood repeater for the node that sent the LIE. In that case the northbound neighbor SHOULD NOT reflood northbound TIEs received from the node that sent the LIE. If the 'you\_are\_flood\_repeater' is absent or if 'you\_are\_flood\_repeater' is set to true, then the northbound neighbor is a flood repeater for the node that sent the LIE and MUST reflood northbound TIEs received from that node.

This specification proposes a simple default algorithm that SHOULD be implemented and used by default on every RIFT node.

- o let  $|NA(Node)$  be the set of Northbound adjacencies of node Node and  $CN(Node)$  be the cardinality of  $|NA(Node)$ ;
- o let  $|SA(Node)$  be the set of Southbound adjacencies of node Node and  $CS(Node)$  be the cardinality of  $|SA(Node)$ ;
- o let  $|P(Node)$  be the set of node Node's parents;
- o let  $|G(Node)$  be the set of node Node's grandparents. Observe that  $|G(Node) = |P(|P(Node))$ ;
- o let N be the child node at level L computing a set of FR;
- o let P be a node at level L+1 and a parent node of N, i.e. bi-directionally reachable over adjacency  $A(N, P)$ ;

- o let G be a grandparent node of N, reachable transitively via a parent P over adjacencies ADJ(N, P) and ADJ(P, G). Observe that N does not have enough information to check bidirectional reachability of ADJ(P, G);
- o let R be a redundancy constant integer; a value of 2 or higher for R is RECOMMENDED;
- o let S be a similarity constant integer; a value in range 0 .. 2 for S is RECOMMENDED, the value of 1 SHOULD be used. Two cardinalities are considered as equivalent if their absolute difference is less than or equal to S, i.e.  $|a-b| \leq S$ .
- o let RND be a 64-bit random number generated by the system once on startup.

The algorithm consists of the following steps:

1. Derive a 64-bits number by XOR'ing 'N's system ID with RND.
2. Derive a 16-bits pseudo-random unsigned integer PR(N) from the resulting 64-bits number by splitting it in 16-bits-long words W1, W2, W3, W4 (where W1 are the least significant 16 bits of the 64-bits number, and W4 are the most significant 16 bits) and then XOR'ing the circularly shifted resulting words together:
  - A.  $(W1 \ll 1) \text{ xor } (W2 \ll 2) \text{ xor } (W3 \ll 3) \text{ xor } (W4 \ll 4)$ ;

where  $\ll$  is the circular shift operator.
3. Sort the parents by decreasing number of northbound adjacencies (using decreasing system id of the parent as tie-breaker):  
 sort  $|P(N)$  by decreasing CN(P), for all P in  $|P(N)$ , as ordered array  $|A(N)$
4. Partition  $|A(N)$  in subarrays  $|A_k(N)$  of parents with equivalent cardinality of northbound adjacencies (in other words with equivalent number of grandparents they can reach):
  - A. set  $k=0$ ; // k is the ID of the subarray
  - B. set  $i=0$ ;
  - C. while  $i < CN(N)$  do
    - i) set  $j=i$ ;
    - ii) while  $i < CN(N)$  and  $CN(|A(N)[j]) - CN(|A(N)[i]) \leq S$

- a. place  $|A(N)[i]$  in  $|A_k(N)$  // abstract action, maybe noop
    - b. set  $i=i+1$ ;
  - iii) /\* At this point  $j$  is the index in  $|A(N)$  of the first member of  $|A_k(N)$  and  $(i-j)$  is  $C_k(N)$  defined as the cardinality of  $|A_k(N)$  \*/
    - set  $k=k+1$ ;
- /\* At this point  $k$  is the total number of subarrays, initialized for the shuffling operation below \*/
- 5. shuffle individually each subarrays  $|A_k(N)$  of cardinality  $C_k(N)$  within  $|A(N)$  using the Durstenfeld variation of Fisher-Yates algorithm that depends on  $N$ 's System ID:
  - A. while  $k > 0$  do
    - i) for  $i$  from  $C_k(N)-1$  to 1 decrementing by 1 do
      - a. set  $j$  to  $PR(N)$  modulo  $i$ ;
      - b. exchange  $|A_k[j]$  and  $|A_k[i]$ ;
    - ii) set  $k=k-1$ ;
- 6. For each grandparent  $G$ , initialize a counter  $c(G)$  with the number of its south-bound adjacencies to elected flood repeaters (which is initially zero):
  - A. for each  $G$  in  $|G(N)$  set  $c(G) = 0$ ;
- 7. Finally keep as FRs only parents that are needed to maintain the number of adjacencies between the FRs and any grandparent  $G$  equal or above the redundancy constant  $R$ :
  - A. for each  $P$  in reshuffled  $|A(N)$ ;
    - i) if there exists an adjacency  $ADJ(P, G)$  in  $|NA(P)$  such that  $c(G) < R$  then
      - a. place  $P$  in FR set;
      - b. for all adjacencies  $ADJ(P, G')$  in  $|NA(P)$  increment  $c(G')$

- B. If any  $c(G)$  is still  $< R$ , it was not possible to elect a set of FRs that covers all grandparents with redundancy  $R$

Additional rules for flooding reduction:

1. The algorithm MUST be re-evaluated by a node on every change of local adjacencies or reception of a parent South TIE with changed adjacencies. A node MAY apply a hysteresis to prevent excessive amount of computation during periods of network instability just like in case of reachability computation.
2. Upon a change of the flood repeater set, a node SHOULD send out LIEs that grant flood repeater status to newly promoted nodes before it sends LIEs that revoke the status to the nodes that have been newly demoted. This is done to prevent transient behavior where the full coverage of grandparents is not guaranteed. Such a condition is sometimes unavoidable in case of lost LIEs but it will correct itself though at possible transient hit in flooding propagation speeds.
3. A node MUST always flood its self-originated TIEs.
4. A node receiving a TIE originated by a node for which it is not a flood repeater SHOULD NOT reflood such TIEs to its neighbors except for rules in Paragraph 6.
5. The indication of flood reduction capability MUST be carried in the node TIEs and MAY be used to optimize the algorithm to account for nodes that will flood regardless.
6. A node generates TIDES as usual but when receiving TIREs or TIDES resulting in requests for a TIE of which the newest received copy came on an adjacency where the node was not flood repeater it SHOULD ignore such requests on first and only first request. Normally, the nodes that received the TIEs as flooding repeaters should satisfy the requesting node and with that no further TIREs for such TIEs will be generated. Otherwise, the next set of TIDES and TIREs MUST lead to flooding independent of the flood repeater status. This solves a very difficult incast problem on nodes restarting with a very wide fanout, especially northbound. To retrieve the full database they often end up processing many in-rushing copies whereas this approach load-balances the incoming database between adjacent nodes and flood repeaters should guarantee that two copies are sent by different nodes to ensure against any losses.

#### 4.2.3.10. Special Considerations

First, due to the distributed, asynchronous nature of ZTP, it can create temporary convergence anomalies where nodes at higher levels of the fabric temporarily see themselves lower than they belong. Since flooding can begin before ZTP is "finished" and in fact must do so given there is no global termination criteria, information may end up in wrong layers. A special clause when changing level takes care of that.

More difficult is a condition where a node (e.g. a leaf) floods a TIE north towards its grandparent, then its parent reboots, in fact partitioning the grandparent from leaf directly and then the leaf itself reboots. That can leave the grandparent holding the "primary copy" of the leaf's TIE. Normally this condition is resolved easily by the leaf re-originating its TIE with a higher sequence number than it sees in northbound TIEs, here however, when the parent comes back it won't be able to obtain leaf's North TIE from the grandparent easily and with that the leaf may not issue the TIE with a higher sequence number that can reach the grandparent for a long time. Flooding procedures are extended to deal with the problem by the means of special clauses that override the database of a lower level with headers of newer TIEs seen in TIEs coming from the north.

#### 4.2.4. Reachability Computation

A node has three possible sources of relevant information for reachability computation. A node knows the full topology south of it from the received North Node TIEs or alternately north of it from the South Node TIEs. A node has the set of prefixes with their associated distances and bandwidths from corresponding prefix TIEs.

To compute prefix reachability, a node runs conceptually a northbound and a southbound SPF. We call that N-SPF and S-SPF denoting the direction in which the computation front is progressing.

Since neither computation can "loop", it is possible to compute non-equal-cost or even k-shortest paths [EPPSTEIN] and "saturate" the fabric to the extent desired but we use simple, familiar SPF algorithms and concepts here as example due to their prevalence in today's routing.

##### 4.2.4.1. Northbound SPF

N-SPF MUST use exclusively northbound and East-West adjacencies in the computing node's node North TIEs (since if the node is a leaf it may not have generated a node South TIE) when starting SPF. Observe that N-SPF is really just a one hop variety since Node South TIEs are



not re-flooded southbound beyond a single level (or East-West) and with that the computation cannot progress beyond adjacent nodes.

Once progressing, we are using the next higher level's node South TIEs to find according adjacencies to verify backlink connectivity. Just as in case of IS-IS or OSPF, two unidirectional links MUST be associated together to confirm bidirectional connectivity. Particular care MUST be paid that the Node TIEs do not only contain the correct system IDs but matching levels as well.

Default route found when crossing an E-W link SHOULD be used IIF

1. the node itself does NOT have any northbound adjacencies AND
2. the adjacent node has one or more northbound adjacencies

This rule forms a "one-hop default route split-horizon" and prevents looping over default routes while allowing for "one-hop protection" of nodes that lost all northbound adjacencies except at Top-of-Fabric where the links are used exclusively to flood topology information in multi-plane designs.

Other south prefixes found when crossing E-W link MAY be used IIF

1. no north neighbors are advertising same or supersuming non-default prefix AND
2. the node does not originate a non-default supersuming prefix itself.

i.e. the E-W link can be used as a gateway of last resort for a specific prefix only. Using south prefixes across E-W link can be beneficial e.g. on automatic de-aggregation in pathological fabric partitioning scenarios.

A detailed example can be found in Section 5.4.

#### 4.2.4.2. Southbound SPF

S-SPF MUST use exclusively the southbound adjacencies in the node South TIEs, i.e. progresses towards nodes at lower levels. Observe that E-W adjacencies are NEVER used in the computation. This enforces the requirement that a packet traversing in a southbound direction must never change its direction.

S-SPF MUST use northbound adjacencies in node North TIEs to verify backlink connectivity by checking for presence of the link beside correct SystemID and level.

#### 4.2.4.3. East-West Forwarding Within a non-ToF Level

Using south prefixes over horizontal links MAY occur if the N-SPF includes East-West adjacencies in computation. It can protect against pathological fabric partitioning cases that leave only paths to destinations that would necessitate multiple changes of forwarding direction between north and south.

#### 4.2.4.4. East-West Links Within ToF Level

E-W ToF links behave in terms of flooding scopes defined in Section 4.2.3.4 like northbound links and MUST be used exclusively for control plane information flooding. Even though a ToF node could be tempted to use those links during southbound SPF and carry traffic over them this MUST NOT be attempted since it may lead in, e.g. anycast cases to routing loops. An implementation MAY try to resolve the looping problem by following on the ring strictly tie-broken shortest-paths only but the details are outside this specification. And even then, the problem of proper capacity provisioning of such links when they become traffic-bearing in case of failures is vexing.

#### 4.2.5. Automatic Disaggregation on Link & Node Failures

##### 4.2.5.1. Positive, Non-transitive Disaggregation

Under normal circumstances, node's South TIEs contain just the adjacencies and a default route. However, if a node detects that its default IP prefix covers one or more prefixes that are reachable through it but not through one or more other nodes at the same level, then it MUST explicitly advertise those prefixes in an South TIE. Otherwise, some percentage of the northbound traffic for those prefixes would be sent to nodes without according reachability, causing it to be black-holed. Even when not black-holing, the resulting forwarding could 'backhaul' packets through the higher level spines, clearly an undesirable condition affecting the blocking probabilities of the fabric.

We refer to the process of advertising additional prefixes southbound as 'positive de-aggregation' or 'positive dis-aggregation'. Such dis-aggregation is non-transitive, i.e. its' effects are always contained to a single level of the fabric only. Naturally, multiple node or link failures can lead to several independent instances of positive dis-aggregation necessary to prevent looping or bow-tying the fabric.

A node determines the set of prefixes needing de-aggregation using the following steps:

1. A DAG computation in the southern direction is performed first, i.e. the North TIEs are used to find all of prefixes it can reach and the set of next-hops in the lower level for each of them. Such a computation can be easily performed on a fat tree by e.g. setting all link costs in the southern direction to 1 and all northern directions to infinity. We term set of those prefixes  $|R$ , and for each prefix,  $r$ , in  $|R$ , we define its set of next-hops to be  $|H(r)$ .
2. The node uses reflected South TIEs to find all nodes at the same level in the same PoD and the set of southbound adjacencies for each. The set of nodes at the same level is termed  $|N$  and for each node,  $n$ , in  $|N$ , we define its set of southbound adjacencies to be  $|A(n)$ .
3. For a given  $r$ , if the intersection of  $|H(r)$  and  $|A(n)$ , for any  $n$ , is null then that prefix  $r$  must be explicitly advertised by the node in an South TIE.
4. Identical set of de-aggregated prefixes is flooded on each of the node's southbound adjacencies. In accordance with the normal flooding rules for an South TIE, a node at the lower level that receives this South TIE SHOULD NOT propagate it south-bound or reflect the disaggregated prefixes back over its adjacencies to nodes at the level from which it was received.

To summarize the above in simplest terms: if a node detects that its default route encompasses prefixes for which one of the other nodes in its level has no possible next-hops in the level below, it has to disaggregate it to prevent black-holing or suboptimal routing through such nodes. Hence a node  $X$  needs to determine if it can reach a different set of south neighbors than other nodes at the same level, which are connected to it via at least one common south neighbor. If it can, then prefix disaggregation may be required. If it can't, then no prefix disaggregation is needed. An example of disaggregation is provided in Section 5.3.

A possible algorithm is described last:

1. Create `partial_neighbors = (empty)`, a set of neighbors with partial connectivity to the node  $X$ 's level from  $X$ 's perspective. Each entry in the set is a south neighbor of  $X$  and a list of nodes of  $X$ .level that can't reach that neighbor.
2. A node  $X$  determines its set of southbound neighbors `X.south_neighbors`.

3. For each South TIE originated from a node Y that X has which is at X.level, if Y.south\_neighbors is not the same as X.south\_neighbors but the nodes share at least one southern neighbor, for each neighbor N in X.south\_neighbors but not in Y.south\_neighbors, add (N, (Y)) to partial\_neighbors if N isn't there or add Y to the list for N.
4. If partial\_neighbors is empty, then node X does not disaggregate any prefixes. If node X is advertising disaggregated prefixes in its South TIE, X SHOULD remove them and re-advertise its according South TIEs.

A node X computes reachability to all nodes below it based upon the received North TIEs first. This results in a set of routes, each categorized by (prefix, path\_distance, next-hop set). Alternately, for clarity in the following procedure, these can be organized by next-hop set as ( (next-hops), {(prefix, path\_distance)}). If partial\_neighbors isn't empty, then the following procedure describes how to identify prefixes to disaggregate.

```
disaggregated_prefixes = { empty }
nodes_same_level = { empty }
for each South TIE
  if (South TIE.level == X.level and
      X shares at least one S-neighbor with X)
    add South TIE.originator to nodes_same_level
  end if
end for

for each next-hop-set NHS
  isolated_nodes = nodes_same_level
  for each NH in NHS
    if NH in partial_neighbors
      isolated_nodes =
        intersection(isolated_nodes,
                    partial_neighbors[NH].nodes)
    end if
  end for

  if isolated_nodes is not empty
    for each prefix using NHS
      add (prefix, distance) to disaggregated_prefixes
    end for
  end if
end for

copy disaggregated_prefixes to X's South TIE
if X's South TIE is different
  schedule South TIE for flooding
end if
```

Figure 15: Computation of Disaggregated Prefixes

Each disaggregated prefix is sent with the according path\_distance. This allows a node to send the same South TIE to each south neighbor. The south neighbor which is connected to that prefix will thus have a shorter path.

Finally, to summarize the less obvious points partially omitted in the algorithms to keep them more tractable:

1. all neighbor relationships MUST perform backlink checks.
2. overload bits as introduced in Section 4.3.1 have to be respected during the computation.

3. all the lower level nodes are flooded the same disaggregated prefixes since we don't want to build an South TIE per node and complicate things unnecessarily. The lower level node that can compute a southbound route to the prefix will prefer it to the disaggregated route anyway based on route preference rules.
4. positively disaggregated prefixes do NOT have to propagate to lower levels. With that the disturbance in terms of new flooding is contained to a single level experiencing failures.
5. disaggregated Prefix South TIEs are not "reflected" by the lower level, i.e. nodes within same level do NOT need to be aware which node computed the need for disaggregation.
6. The fabric is still supporting maximum load balancing properties while not trying to send traffic northbound unless necessary.

In case positive disaggregation is triggered and due to the very stable but un-synchronized nature of the algorithm the nodes may issue the necessary disaggregated prefixes at different points in time. This can lead for a short time to an "incast" behavior where the first advertising router based on the nature of longest prefix match will attract all the traffic. An implementation MAY hence choose different strategies to address this behavior if needed.

To close this section it is worth to observe that in a single plane ToF this disaggregation prevents blackholing up to  $(K\_LEAF * P)$  link failures in terms of Section 4.1.2 or in other terms, it takes at minimum that many link failures to partition the ToF into multiple planes.

#### 4.2.5.2. Negative, Transitive Disaggregation for Fallen Leaves

As explained in Section 4.1.3 failures in multi-plane Top-of-Fabric or more than  $(K\_LEAF * P)$  links failing in single plane design can generate fallen leaves. Such scenario cannot be addressed by positive disaggregation only and needs a further mechanism.

##### 4.2.5.2.1. Cabling of Multiple Top-of-Fabric Planes

Let us return in this section to designs with multiple planes as shown in Figure 3. Figure 16 highlights how the ToF is cabled in case of two planes by the means of dual-rings to distribute all the North TIEs within both planes. For people familiar with traditional link-state routing protocols ToF level can be considered equivalent to area 0 in OSPF or level-2 in ISIS which need to be "connected" as well for the protocol to operate correctly.

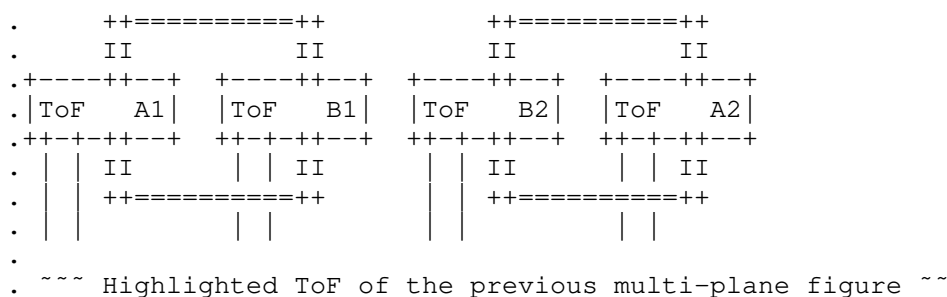


Figure 16: Topologically Connected Planes

As described in Section 4.1.3 failures in multi-plane fabrics can lead to blackholes which normal positive disaggregation cannot fix. The mechanism of negative, transitive disaggregation incorporated in RIFT provides the according solution.

#### 4.2.5.2.2. Transitive Advertisement of Negative Disaggregates

A ToF node that discovers that it cannot reach a fallen leaf disaggregates all the prefixes of such leaves. It uses for that purpose negative prefix South TIEs that are, as usual, flooded southwards with the scope defined in Section 4.2.3.4.

Transitively, a node explicitly loses connectivity to a prefix when none of its children advertises it and when the prefix is negatively disaggregated by all of its parents. When that happens, the node originates the negative prefix further down south. Since the mechanism applies recursively south the negative prefix may propagate transitively all the way down to the leaf. This is necessary since leaves connected to multiple planes by means of disjoint paths may have to choose the correct plane already at the very bottom of the fabric to make sure that they don't send traffic towards another leaf using a plane where it is "fallen" at which in point a blackhole is unavoidable.

When the connectivity is restored, a node that disaggregated a prefix withdraws the negative disaggregation by the usual mechanism of re-advertising TIEs omitting the negative prefix.

#### 4.2.5.2.3. Computation of Negative Disaggregates

The document omitted so far the description of the computation necessary to generate the correct set of negative prefixes. Negative prefixes can in fact be advertised due to two different triggers. We describe them consecutively.

The first origination reason is a computation that uses all the node North TIEs to build the set of all reachable nodes by reachability computation over the complete graph and including ToF links. The computation uses the node itself as root. This is compared with the result of the normal southbound SPF as described in Section 4.2.4.2. The difference are the fallen leaves and all their attached prefixes are advertised as negative prefixes southbound if the node does not see the prefix being reachable within southbound SPF.

The second mechanism hinges on the understanding how the negative prefixes are used within the computation as described in Figure 17. When attaching the negative prefixes at certain point in time the negative prefix may find itself with all the viable nodes from the shorter match nexthop being pruned. In other words, all its northbound neighbors provided a negative prefix advertisement. This is the trigger to advertise this negative prefix transitively south and normally caused by the node being in a plane where the prefix belongs to a fabric leaf that has "fallen" in this plane. Obviously, when one of the northbound switches withdraws its negative advertisement, the node has to withdraw its transitively provided negative prefix as well.

#### 4.2.6. Attaching Prefixes

After SPF is run, it is necessary to attach the resulting reachability information in form of prefixes. For S-SPF, prefixes from an North TIE are attached to the originating node with that node's next-hop set and a distance equal to the prefix's cost plus the node's minimized path distance. The RIFT route database, a set of (prefix, prefix-type, attributes, path\_distance, next-hop set), accumulates these results.

In case of N-SPF prefixes from each South TIE need to also be added to the RIFT route database. The N-SPF is really just a stub so the computing node needs simply to determine, for each prefix in an South TIE that originated from adjacent node, what next-hops to use to reach that node. Since there may be parallel links, the next-hops to use can be a set; presence of the computing node in the associated Node South TIE is sufficient to verify that at least one link has bidirectional connectivity. The set of minimum cost next-hops from the computing node X to the originating adjacent node is determined.

Each prefix has its cost adjusted before being added into the RIFT route database. The cost of the prefix is set to the cost received plus the cost of the minimum distance next-hop to that neighbor while taking into account its attributes such as mobility per Section 4.3.3. Then each prefix can be added into the RIFT route database with the next-hop set; ties are broken based upon type first



and then distance and further on 'PrefixAttributes' and only the best combination is used for forwarding. RIFT route preferences are normalized by the according Thrift [thrift] model type.

An example implementation for node X follows:

```

for each South TIE
  if South TIE.level > X.level
    next_hop_set = set of minimum cost links to the
                  South TIE.originator
    next_hop_cost = minimum cost link to
                  South TIE.originator
  end if
  for each prefix P in the South TIE
    P.cost = P.cost + next_hop_cost
    if P not in route_database:
      add (P, P.cost, P.type,
          P.attributes, next_hop_set) to route_database
    end if
    if (P in route_database):
      if route_database[P].cost > P.cost or
         route_database[P].type > P.type:
        update route_database[P] with (P, P.type, P.cost,
                                       P.attributes,
                                       next_hop_set)
      else if route_database[P].cost == P.cost and
              route_database[P].type == P.type:
        update route_database[P] with (P, P.type,
                                       P.cost, P.attributes,
                                       merge(next_hop_set, route_database[P].next_hop_set))
      else
        // Not preferred route so ignore
      end if
    end if
  end for
end for

```

Figure 17: Adding Routes from South TIE Positive and Negative Prefixes

After the positive prefixes are attached and tie-broken, negative prefixes are attached and used in case of northbound computation, ideally from the shortest length to the longest. The nexthop adjacencies for a negative prefix are inherited from the longest positive prefix that aggregates it, and subsequently adjacencies to nodes that advertised negative for this prefix are removed.

The rule of inheritance MUST be maintained when the nexthop list for a prefix is modified, as the modification may affect the entries for matching negative prefixes of immediate longer prefix length. For instance, if a nexthop is added, then by inheritance it must be added to all the negative routes of immediate longer prefixes length unless it is pruned due to a negative advertisement for the same next hop. Similarly, if a nexthop is deleted for a given prefix, then it is deleted for all the immediately aggregated negative routes. This will recurse in the case of nested negative prefix aggregations.

The rule of inheritance must also be maintained when a new prefix of intermediate length is inserted, or when the immediately aggregating prefix is deleted from the routing table, making an even shorter aggregating prefix the one from which the negative routes now inherit their adjacencies. As the aggregating prefix changes, all the negative routes must be recomputed, and then again the process may recurse in case of nested negative prefix aggregations.

Although these operations can be computationally expensive, the overall load on devices in the network is low because these computations are not run very often, as positive route advertisements are always preferred over negative ones. This prevents recursion in most cases because positive reachability information never inherits next hops.

To make the negative disaggregation less abstract and provide an example let us consider a ToP node T1 with 4 ToF parents S1..S4 as represented in Figure 18:

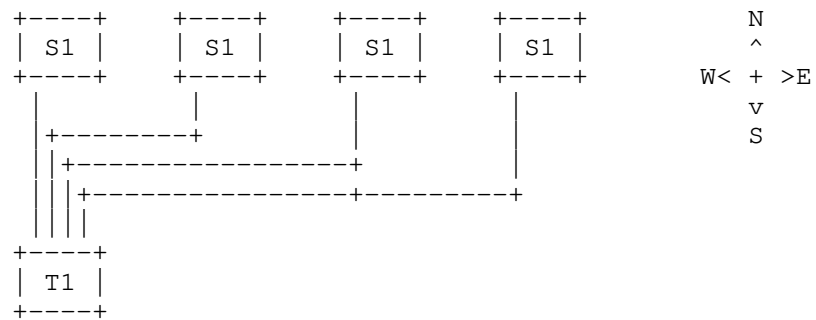


Figure 18: A ToP Node with 4 Parents

If all ToF nodes can reach all the prefixes in the network; with RIFT, they will normally advertise a default route south. An abstract Routing Information Base (RIB), more commonly known as a

routing table, stores all types of maintained routes including the negative ones and "tie-breaks" for the best one, whereas an abstract Forwarding table (FIB) retains only the ultimately computed "positive" routing instructions. In T1, those tables would look as illustrated in Figure 19:

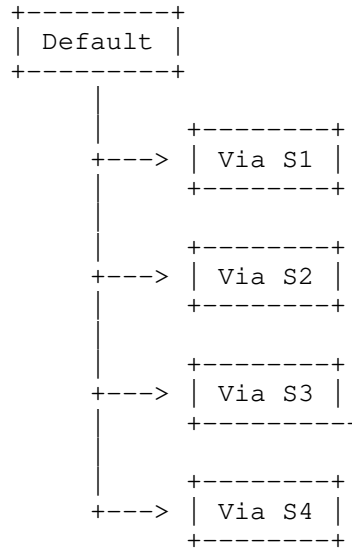


Figure 19: Abstract RIB

In case T1 receives a negative advertisement for prefix 2001:db8::/32 from S1 a negative route is stored in the RIB (indicated by a ~ sign), while the more specific routes to the complementing ToF nodes are installed in FIB. RIB and FIB in T1 now look as illustrated in Figure 20 and Figure 21, respectively:

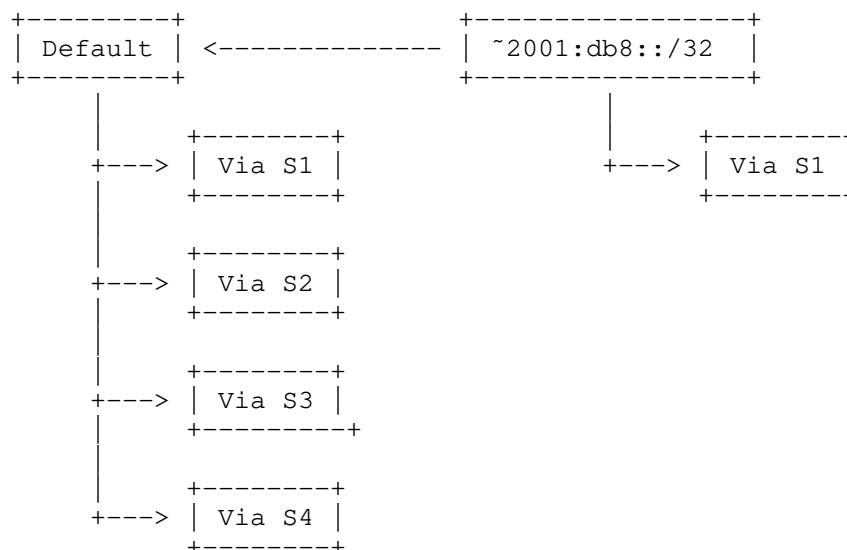


Figure 20: Abstract RIB after Negative 2001:db8::/32 from S1

The negative 2001:db8::/32 prefix entry inherits from ::/0, so the positive more specific routes are the complements to S1 in the set of next-hops for the default route. That entry is composed of S2, S3, and S4, or, in other words, it uses all entries the the default route with a "hole punched" for S1 into them. These are the next hops that are still available to reach 2001:db8::/32, now that S1 advertised that it will not forward 2001:db8::/32 anymore. Ultimately, those resulting next-hops are installed in FIB for the more specific route to 2001:db8::/32 as illustrated below:

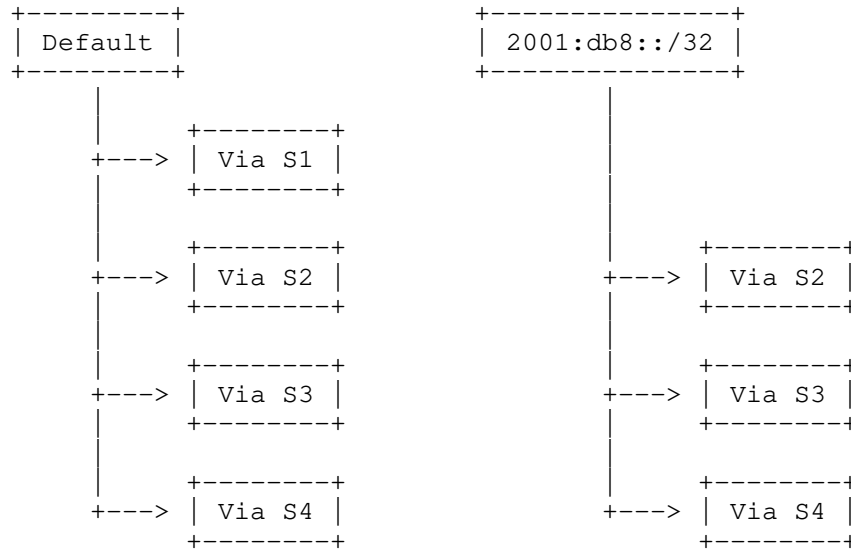


Figure 21: Abstract FIB after Negative 2001:db8::/32 from S1

To illustrate matters further let us consider T1 receiving a negative advertisement for prefix 2001:db8:1::/48 from S2, which is stored in RIB again. After the update, the RIB in T1 is illustrated in Figure 22:

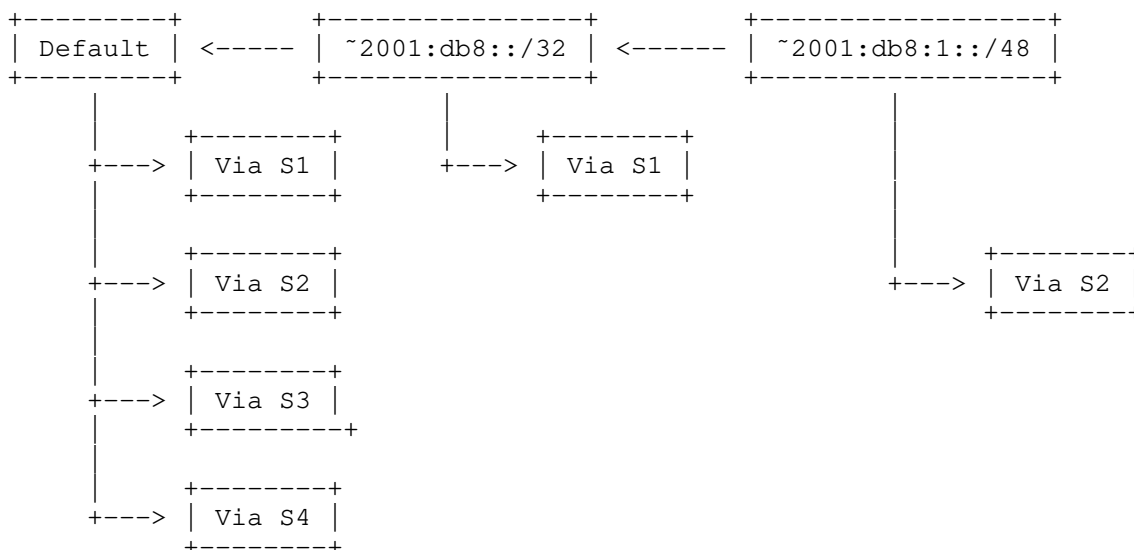


Figure 22: Abstract RIB after Negative 2001:db8:1::/48 from S2

Negative 2001:db8:1::/48 inherits from 2001:db8::/32 now, so the positive more specific routes are the complements to S2 in the set of next hops for 2001:db8::/32, which are S3 and S4, or, in other words, all entries of the parent with the negative holes "punched in" again. After the update, the FIB in T1 shows as illustrated in Figure 23:







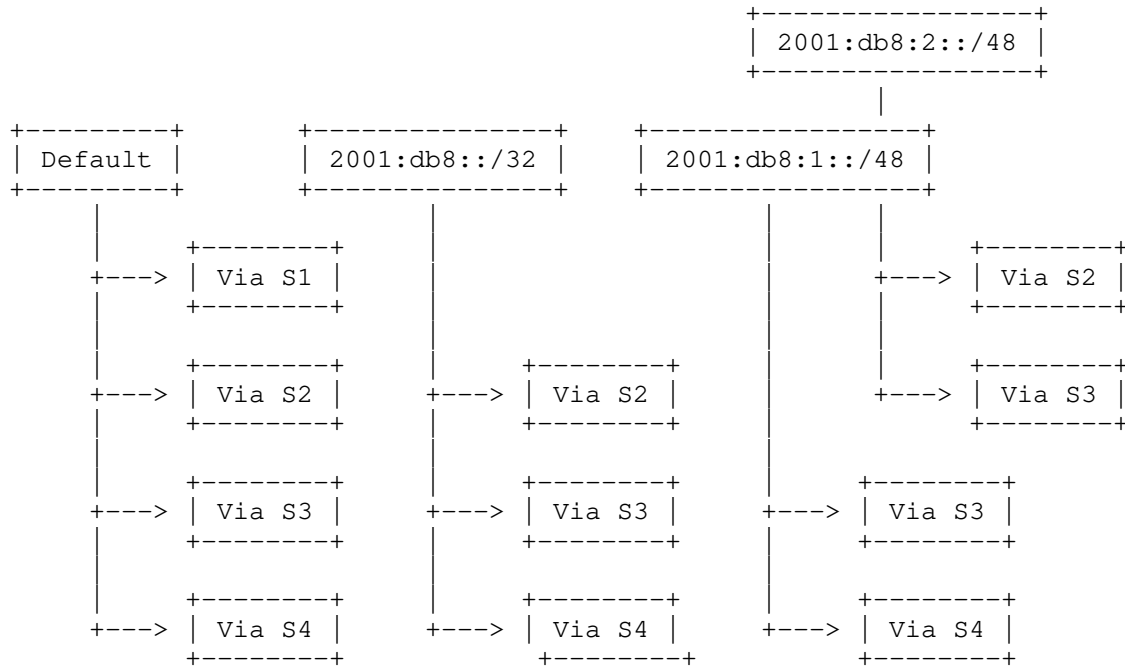


Figure 25: Abstract FIB after Negative `2001:db8:2::/48` from S4

4.2.7. Optional Zero Touch Provisioning (ZTP)

Each RIFT node can operate in zero touch provisioning (ZTP) mode, i.e. it has no configuration (unless it is a ToF or it is configured to operate in the overall topology as leaf and/or support leaf-2-leaf procedures) and it will fully configure itself after being attached to the topology. Configured nodes and nodes operating in ZTP can be mixed and will form a valid topology if achievable.

The derivation of the level of each node happens based on offers received from its neighbors whereas each node (with possibly exceptions of configured leaves) tries to attach at the highest possible point in the fabric. This guarantees that even if the diffusion front reaches a node from "below" faster than from "above", it will greedily abandon already negotiated level derived from nodes topologically below it and properly peers with nodes above.

The fabric is very consciously numbered from the top to allow for PoDs of different heights and minimize number of provisioning

necessary, in this case just a TOP\_OF\_FABRIC flag on every node at the top of the fabric.

This section describes the necessary concepts and procedures for ZTP operation.

#### 4.2.7.1. Terminology

The interdependencies between the different flags and the configured level can be somewhat vexing at first and it may take multiple reads of the glossary to comprehend them.

**Automatic Level Derivation:** Procedures which allow nodes without level configured to derive it automatically. Only applied if CONFIGURED\_LEVEL is undefined.

**UNDEFINED\_LEVEL:** A "null" value that indicates that the level has not been determined and has not been configured. Schemas normally indicate that by a missing optional value without an available defined default.

**LEAF\_ONLY:** An optional configuration flag that can be configured on a node to make sure it never leaves the "bottom of the hierarchy". TOP\_OF\_FABRIC flag and CONFIGURED\_LEVEL cannot be defined at the same time as this flag. It implies CONFIGURED\_LEVEL value of 0.

**TOP\_OF\_FABRIC flag:** Configuration flag that MUST be provided to all Top-of-Fabric nodes. LEAF\_FLAG and CONFIGURED\_LEVEL cannot be defined at the same time as this flag. It implies a CONFIGURED\_LEVEL value. In fact, it is basically a shortcut for configuring same level at all Top-of-Fabric nodes which is unavoidable since an initial 'seed' is needed for other ZTP nodes to derive their level in the topology. The flag plays an important role in fabrics with multiple planes to enable successful negative disaggregation (Section 4.2.5.2).

**CONFIGURED\_LEVEL:** A level value provided manually. When this is defined (i.e. it is not an UNDEFINED\_LEVEL) the node is not participating in ZTP. TOP\_OF\_FABRIC flag is ignored when this value is defined. LEAF\_ONLY can be set only if this value is undefined or set to 0.

**DERIVED\_LEVEL:** Level value computed via automatic level derivation when CONFIGURED\_LEVEL is equal to UNDEFINED\_LEVEL.

**LEAF\_2\_LEAF:** An optional flag that can be configured on a node to make sure it supports procedures defined in Section 4.3.8. In a strict sense it is a capability that implies LEAF\_ONLY and the

according restrictions. TOP\_OF\_FABRIC flag is ignored when set at the same time as this flag.

**LEVEL\_VALUE:** In ZTP case the original definition of "level" in Section 3.1 is both extended and relaxed. First, level is defined now as LEVEL\_VALUE and is the first defined value of CONFIGURED\_LEVEL followed by DERIVED\_LEVEL. Second, it is possible for nodes to be more than one level apart to form adjacencies if any of the nodes is at least LEAF\_ONLY.

**Valid Offered Level (VOL):** A neighbor's level received on a valid LIE (i.e. passing all checks for adjacency formation while disregarding all clauses involving level values) persisting for the duration of the holdtime interval on the LIE. Observe that offers from nodes offering level value of 0 do not constitute VOLs (since no valid DERIVED\_LEVEL can be obtained from those and consequently 'not\_a\_ztp\_offer' MUST be ignored). Offers from LIEs with 'not\_a\_ztp\_offer' being true are not VOLs either. If a node maintains parallel adjacencies to the neighbor, VOL on each adjacency is considered as equivalent, i.e. the newest VOL from any such adjacency updates the VOL received from the same node.

**Highest Available Level (HAL):** Highest defined level value seen from all VOLs received.

**Highest Available Level Systems (HALS):** Set of nodes offering HAL VOLs.

**Highest Adjacency Three Way (HAT):** Highest neighbor level of all the formed three way adjacencies for the node.

#### 4.2.7.2. Automatic SystemID Selection

RIFT nodes require a 64 bit SystemID which SHOULD be derived as EUI-64 MA-L derive according to [EUI64]. The organizationally governed portion of this ID (24 bits) can be used to generate multiple IDs if required to indicate more than one RIFT instance."

As matter of operational concern, the router MUST ensure that such identifier is not changing very frequently (or at least not without sending all its TIEs with fairly short lifetimes) since otherwise the network may be left with large amounts of stale TIEs in other nodes (though this is not necessarily a serious problem if the procedures described in Section 7 are implemented).

4.2.7.3. Generic Fabric Example

ZTP forces us to think about miscabled or unusually cabled fabric and how such a topology can be forced into a "lattice" structure which a fabric represents (with further restrictions). Let us consider a necessary and sufficient physical cabling in Figure 26. We assume all nodes being in the same PoD.

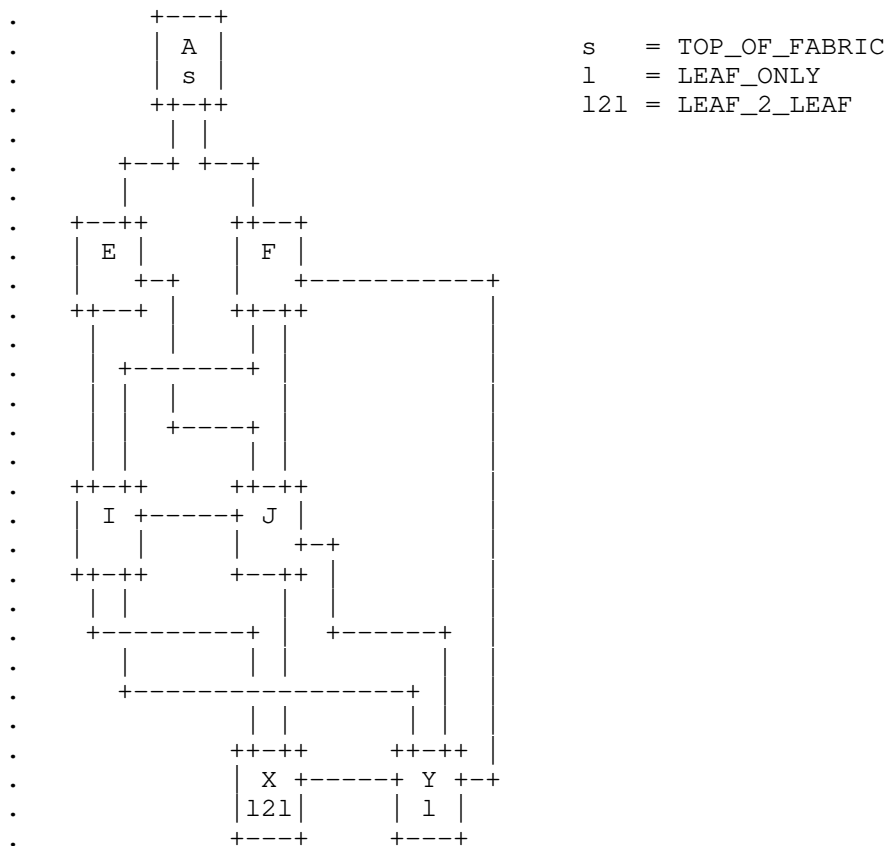


Figure 26: Generic ZTP Cabling Considerations

First, we must anchor the "top" of the cabling and that's what the TOP\_OF\_FABRIC flag at node A is for. Then things look smooth until we have to decide whether node Y is at the same level as I, J (and as consequence, X is south of it) or at the same level as X. This is unresolvable here until we "nail down the bottom" of the topology. To achieve that we choose to use in this example the leaf flags in X and Y. In case where Y would not have a leaf flag it will try to

elect highest level offered and end up being in same level as I and J.

#### 4.2.7.4. Level Determination Procedure

A node starting up with UNDEFINED\_VALUE (i.e. without a CONFIGURED\_LEVEL or any leaf or TOP\_OF\_FABRIC flag) MUST follow those additional procedures:

1. It advertises its LEVEL\_VALUE on all LIEs (observe that this can be UNDEFINED\_LEVEL which in terms of the schema is simply an omitted optional value).
2. It computes HAL as numerically highest available level in all VOLs.
3. It chooses then  $\text{MAX}(\text{HAL}-1, 0)$  as its DERIVED\_LEVEL. The node then starts to advertise this derived level.
4. A node that lost all adjacencies with HAL value MUST hold down computation of new DERIVED\_LEVEL for a short period of time unless it has no VOLs from southbound adjacencies. After the holddown expired, it MUST discard all received offers, recompute DERIVED\_LEVEL and announce it to all neighbors.
5. A node MUST reset any adjacency that has changed the level it is offering and is in three-way state.
6. A node that changed its defined level value MUST readvertise its own TIEs (since the new 'PacketHeader' will contain a different level than before). Sequence number of each TIE MUST be increased.
7. After a level has been derived the node MUST set the 'not\_a\_ztp\_offer' on LIEs towards all systems offering a VOL for HAL.
8. A node that changed its level SHOULD flush from its link state database TIEs of all other nodes, otherwise stale information may persist on "direction reversal", i.e. nodes that seemed south are now north or east-west. This will not prevent the correct operation of the protocol but could be slightly confusing operationally.

A node starting with LEVEL\_VALUE being 0 (i.e. it assumes a leaf function by being configured with the appropriate flags or has a CONFIGURED\_LEVEL of 0) MUST follow those additional procedures:

1. It computes HAT per procedures above but does NOT use it to compute DERIVED\_LEVEL. HAT is used to limit adjacency formation per Section 4.2.2.

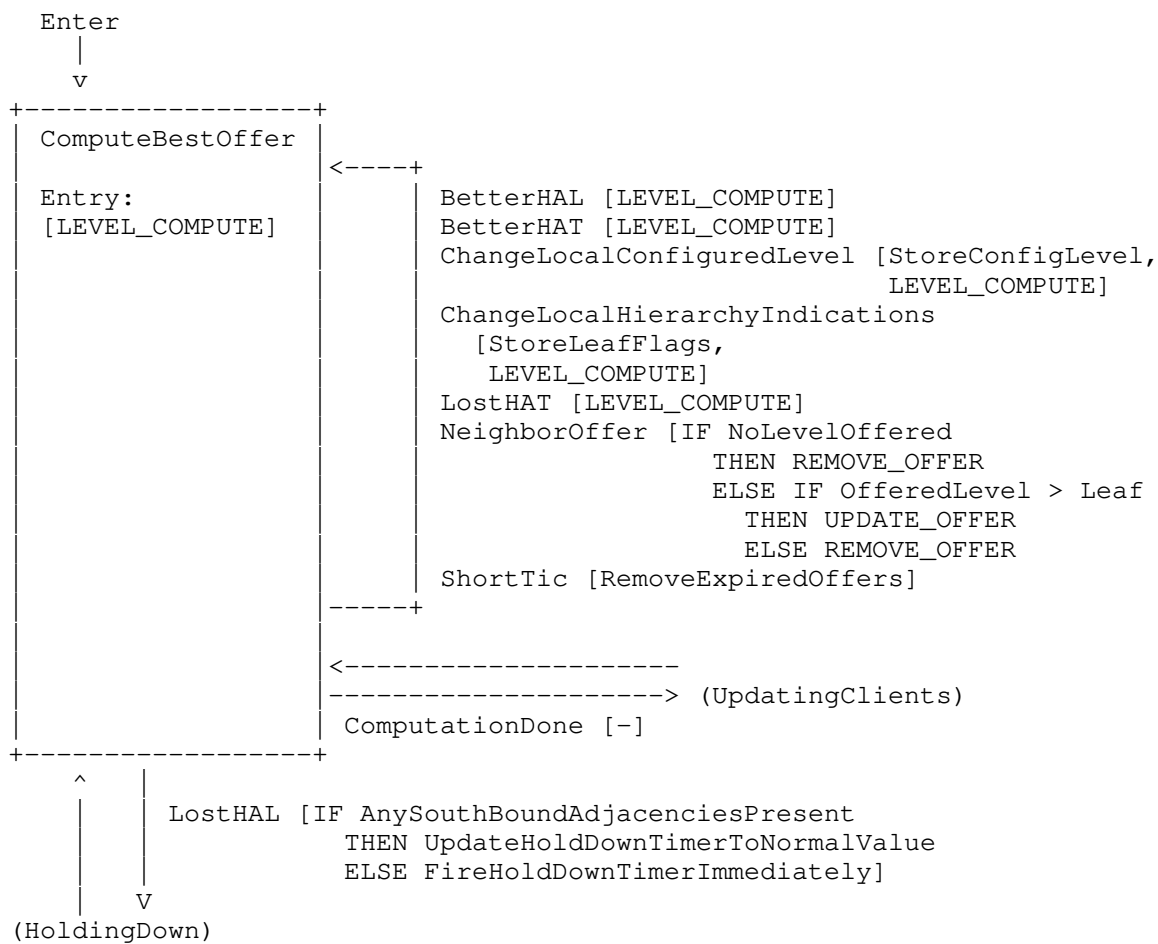
It MAY also follow modified procedures:

1. It may pick a different strategy to choose VOL, e.g. use the VOL value with highest number of VOLs. Such strategies are only possible since the node always remains "at the bottom of the fabric" while another layer could "invert" the fabric by picking its preferred VOL in a different fashion than always trying to achieve the highest viable level.

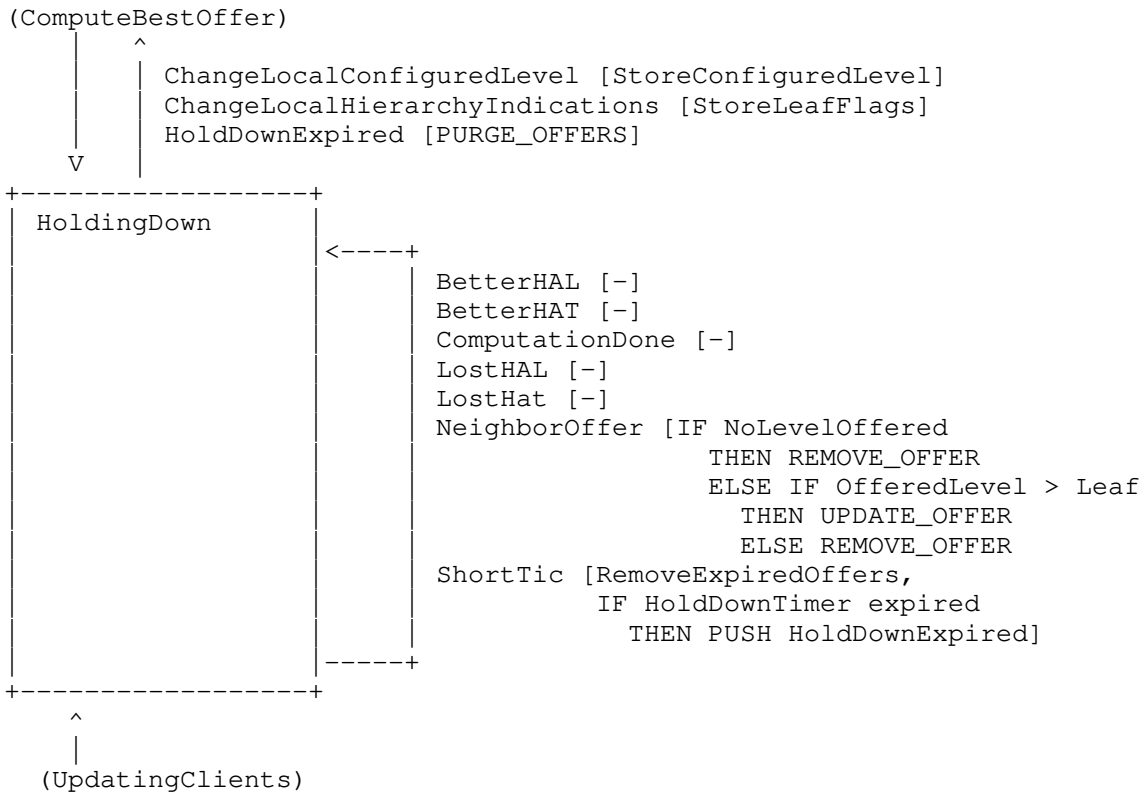
#### 4.2.7.5. ZTP FSM

This section specifies the precise, normative ZTP FSM and can be omitted unless the reader is pursuing an implementation of the protocol.

Initial state is ComputeBestOffer.

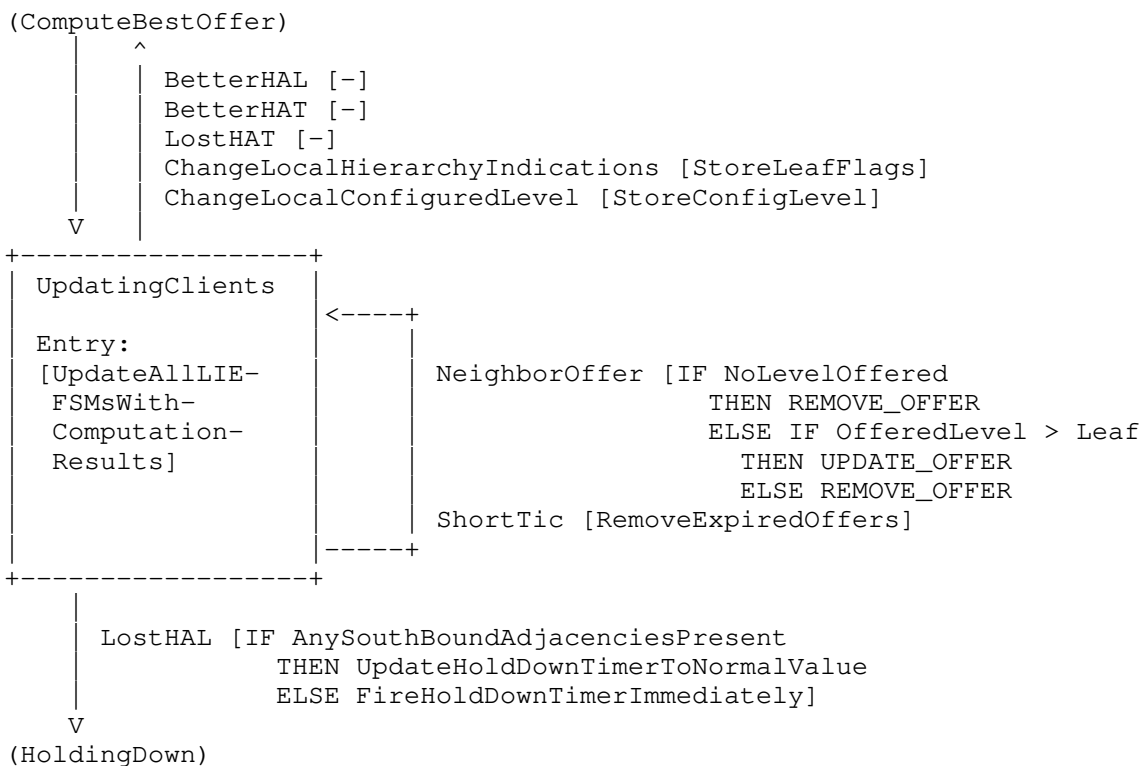


ZTP FSM FSM



ZTP FSM FSM (continued)





ZTP FSM FSM (continued)

Events

- o ChangeLocalHierarchyIndications: node locally configured with new leaf flags
- o ChangeLocalConfiguredLevel: node locally configured with a defined level
- o NeighborOffer: a new neighbor offer with optional level and neighbor state
- o BetterHAL: better HAL computed internally
- o BetterHAT: better HAT computed internally
- o LostHAL: lost last HAL in computation
- o LostHAT: lost HAT in computation

- o ComputationDone: computation performed
- o HoldDownExpired: holddown expired
- o ShortTic: one second timer tick, to be ignored if transition does not exist

#### Actions

on ShortTic in HoldingDown finishes in HoldingDown: remove expired offers and if holddown timer expired PUSH\_EVENT HoldDownExpired

on ShortTic in ComputeBestOffer finishes in ComputeBestOffer: remove expired offers

on HoldDownExpired in HoldingDown finishes in ComputeBestOffer: PURGE\_OFFERS

on ChangeLocalConfiguredLevel in HoldingDown finishes in ComputeBestOffer: store configured level

on ShortTic in UpdatingClients finishes in UpdatingClients: remove expired offers

on BetterHAT in ComputeBestOffer finishes in ComputeBestOffer: LEVEL\_COMPUTE

on BetterHAL in HoldingDown finishes in HoldingDown: no action

on ChangeLocalHierarchyIndications in HoldingDown finishes in ComputeBestOffer: store leaf flags

on BetterHAT in UpdatingClients finishes in ComputeBestOffer: no action

on BetterHAL in UpdatingClients finishes in ComputeBestOffer: no action

on ChangeLocalHierarchyIndications in UpdatingClients finishes in ComputeBestOffer: store leaf flags

on LostHAL in HoldingDown finishes in HoldingDown:

on LostHAT in ComputeBestOffer finishes in ComputeBestOffer: LEVEL\_COMPUTE

on LostHAT in HoldingDown finishes in HoldingDown: no action

```
on BetterHAT in HoldingDown finishes in HoldingDown: no action
on NeighborOffer in UpdatingClients finishes in UpdatingClients:
    if no level offered then REMOVE_OFFER
else
    if offered level > leaf then UPDATE_OFFER
    else REMOVE_OFFER
on LostHAL in ComputeBestOffer finishes in HoldingDown: if any
southbound adjacencies present then update holddown timer to
normal duration else fire holddown timer immediately
on LostHAL in UpdatingClients finishes in HoldingDown: if any
southbound adjacencies present then update holddown timer to
normal duration else fire holddown timer immediately
on ComputationDone in ComputeBestOffer finishes in
UpdatingClients: no action
on LostHAT in UpdatingClients finishes in ComputeBestOffer: no
action
on ComputationDone in HoldingDown finishes in HoldingDown:
on ChangeLocalConfiguredLevel in ComputeBestOffer finishes in
ComputeBestOffer: store configured level and LEVEL_COMPUTE
on ChangeLocalConfiguredLevel in UpdatingClients finishes in
ComputeBestOffer: store configured level
on NeighborOffer in ComputeBestOffer finishes in ComputeBestOffer:
    if no level offered then REMOVE_OFFER
else
    if offered level > leaf then UPDATE_OFFER
    else REMOVE_OFFER
on NeighborOffer in HoldingDown finishes in HoldingDown:
    if no level offered then REMOVE_OFFER
```

else

    if offered level > leaf then UPDATE\_OFFER

    else REMOVE\_OFFER

on ChangeLocalHierarchyIndications in ComputeBestOffer finishes in ComputeBestOffer: store leaf flags and LEVEL\_COMPUTE

on BetterHAL in ComputeBestOffer finishes in ComputeBestOffer: LEVEL\_COMPUTE

on Entry into UpdatingClients: update all LIE FSMs with computation results

on Entry into ComputeBestOffer: LEVEL\_COMPUTE

Following words are used for well known procedures:

1. PUSH Event: pushes an event to be executed by the FSM upon exit of this action
2. COMPARE\_OFFERS: checks whether based on current offers and held last results the events BetterHAL/LostHAL/BetterHAT/LostHAT are necessary and returns them
3. UPDATE\_OFFER: store current offer with adjacency holdtime as lifetime and COMPARE\_OFFERS, then PUSH according events
4. LEVEL\_COMPUTE: compute best offered or configured level and HAL/HAT, if anything changed PUSH ComputationDone
5. REMOVE\_OFFER: remove the according offer and COMPARE\_OFFERS, PUSH according events
6. PURGE\_OFFERS: REMOVE\_OFFER for all held offers, COMPARE OFFERS, PUSH according events

#### 4.2.7.6. Resulting Topologies

The procedures defined in Section 4.2.7.4 will lead to the RIFT topology and levels depicted in Figure 27.





### 4.3. Further Mechanisms

#### 4.3.1. Overload Bit

The overload bit **MUST** be respected by all necessary SPF computations. A node with the overload bit set **SHOULD** advertise all locally hosted prefixes both northbound and southbound, all other southbound prefixes **SHOULD NOT** be advertised.

Leaf nodes **SHOULD** set the overload bit on all originated Node TIEs. If spine nodes were to forward traffic not intended for the local node, the leaf node would not be able to prevent routing/forwarding loops as it does not have the necessary topology information to do so.

#### 4.3.2. Optimized Route Computation on Leaves

Leaf nodes only have visibility to directly connected nodes and therefore are not required to run "full" SPF computations. Instead, prefixes from neighboring nodes can be gathered to run a "partial" SPF computation in order to build the routing table.

Leaf nodes **SHOULD** only hold their own N-TIEs, and in cases of L2L implementations, the N-TIEs of their East/West neighbors. Leaf nodes **MUST** hold all S-TIEs from their neighbors.

Normally, a full network graph is created based on local N-TIEs and remote S-TIEs that it receives from neighbors, at which time, necessary SPF computations are performed. Instead, leaf nodes can simply compute the minimum cost and next-hop set of each leaf neighbor by examining its local adjacencies. Associated N-TIEs are used to determine bi-directionality and derive the next-hop set. Cost is then derived from the minimum cost of the local adjacency to the neighbor and the prefix cost.

Leaf nodes would then attach necessary prefixes as described in Section 4.2.6.

#### 4.3.3. Mobility

The RIFT control plane **MUST** maintain the real time status of every prefix, to which port it is attached, and to which leaf node that port belongs. This is still true in cases of IP mobility where the point of attachment may change several times a second.

There are two classic approaches to explicitly maintain this information:

timestamp: With this method, the infrastructure SHOULD record the precise time at which the movement is observed. One key advantage of this technique is that it has no dependency on the mobile device. One drawback is that the infrastructure MUST be precisely synchronized in order to be able to compare timestamps as the points of attachment change. This could be accomplished by utilizing Precision Time Protocol (PTP) IEEE Std. 1588 [IEEEstd1588] or 802.1AS [IEEEstd8021AS] which is designed for bridged LANs. Both the precision of the synchronization protocol and the resolution of the timestamp must beat the highest possible roaming time on the fabric. Another drawback is that the presence of a mobile device may only be observed asynchronously, such as when it starts using an IP protocol like ARP [RFC0826], IPv6 Neighbor Discovery [RFC4861], IPv6 Stateless Address Configuration [RFC4862], DHCP [RFC2131], or DHCPv6 [RFC8415].

sequence counter: With this method, a mobile device notifies its point of attachment on arrival with a sequence counter that is incremented upon each movement. On the positive side, this method does not have a dependency on a precise sense of time, since the sequence of movements is kept in order by the mobile device. The disadvantage of this approach is the lack of support for protocols that may be used by the mobile device to register its presence to the leaf node with the capability to provide a sequence counter. Well-known issues with sequence counters such as wrapping and comparison rules MUST be addressed properly. Sequence numbers MUST be compared by a single homogenous source to make operation feasible. Sequence number comparison from multiple heterogeneous sources would be extremely difficult to implement.

RIFT supports a hybrid approach by using an optional 'PrefixSequenceType' attribute (that we also call a 'monotonic clock') that consists of a timestamp and optional sequence number field. When this attribute is present (observe that per data schema the attribute itself is optional but in case it is included the 'timestamp' field is required):

- o The leaf node MAY advertise a timestamp of the latest sighting of a prefix, e.g., by snooping IP protocols or the node using the time at which it advertised the prefix. RIFT transports the timestamp within the desired prefix North TIEs as 802.1AS timestamp.
- o RIFT MAY interoperate with "Registration Extensions for 6LoWPAN Neighbor Discovery" [RFC8505], which provides a method for registering a prefix with a sequence number called a Transaction ID (TID). In such cases, RIFT SHOULD transport the derived TID without modification.



- o RIFT also defines an abstract negative clock (ASNC) (also called an 'undefined' clock). ASNC MUST be considered older than any other defined clock. By default, when a node receives a prefix North TIE that does not contain a 'PrefixSequenceType' attribute, it MUST interpret the absence as ASNC.
- o Any prefix present on the fabric in multiple nodes that has the 'same' clock is considered as anycast.
- o RIFT specification assumes that all nodes are being synchronized to at least 200 milliseconds of precision. This is achievable through the use of NTP [RFC5905]. An implementation MAY provide a way to reconfigure a domain to a different value, we call this variable MAXIMUM\_CLOCK\_DELTA.

#### 4.3.3.1. Clock Comparison

All monotonic clock values MUST be compared to each other using the following rules:

1. ASNC is older than any other value except ASNC AND
2. Clock with timestamp differing by more than MAXIMUM\_CLOCK\_DELTA are comparable by using the timestamps only AND
3. Clocks with timestamps differing by less than MAXIMUM\_CLOCK\_DELTA are comparable by using their TIDs only AND
4. An undefined TID is always older than any other TID AND
5. TIDs are compared using rules of [RFC8505].

#### 4.3.3.2. Interaction between Time Stamps and Sequence Counters

For attachment changes that occur less frequently (e.g. once per second), the timestamp that the RIFT infrastructure captures should be enough to determine the most current discovery. If the point of attachment changes faster than the maximum drift of the time stamping mechanism (i.e. MAXIMUM\_CLOCK\_DELTA), then a sequence number SHOULD be used to enable necessary precision to determine currency.

The sequence counter in [RFC8505] is encoded as one octet and wraps around using Appendix A.

Within the resolution of MAXIMUM\_CLOCK\_DELTA, sequence counter values captured during 2 sequential iterations of the same timestamp SHOULD be comparable. This means that with default values, a node may move up to 127 times in a 200 millisecond period and the clocks will

remain comparable. This allows the RIFT infrastructure to explicitly assert the most up-to-date advertisement.

#### 4.3.3.3. Anycast vs. Unicast

A unicast prefix can be attached to at most one leaf, whereas an anycast prefix may be reachable via more than one leaf.

If a monotonic clock attribute is provided on the prefix, then the prefix with the 'newest' clock value is strictly preferred. An anycast prefix does not carry a clock or all clock attributes MUST be the same under the rules of Section 4.3.3.1.

Observe that it is important that in mobility events the leaf is re-flooding as quickly as possible the absence of the prefix that moved away.

Observe further that without support for [RFC8505] movements on the fabric within intervals smaller than 100msec will be seen as anycast.

#### 4.3.3.4. Overlays and Signaling

RIFT is agnostic to any overlay technologies and their associated control and transports that run on top of it (e.g. VXLAN). It is expected that leaf nodes and possibly Top-of-Fabric nodes can perform necessary data plane encapsulation.

In the context of mobility, overlays provide another possible solution to avoid injecting mobile prefixes into the fabric as well as improving scalability of the deployment. It makes sense to consider overlays for mobility solutions in IP fabrics. As an example, a mobility protocol such as LISP [RFC6830] may inform the ingress leaf of the location of the egress leaf in real time.

Another possibility is to consider that mobility as an underlay service and support it in RIFT to an extent. The load on the fabric augments with the amount of mobility obviously since a move forces flooding and computation on all nodes in the scope of the move so tunneling from leaf to the Top-of-Fabric may be desired to speed up convergence times.

#### 4.3.4. Key/Value Store

##### 4.3.4.1. Southbound

RIFT supports the southbound distribution of key-value pairs that can be used to distribute information to facilitate higher levels of functionality (e.g. distribution of configuration information). KV

South TIEs may arrive from multiple nodes and therefore MUST execute the following tie-breaking rules for each key:

1. Only KV TIEs received from nodes to which a bi-directional adjacency exists MUST be considered.
2. For each valid KV South TIEs that contains the same key, the value within the South TIE with the highest level will be preferred. If the levels are identical, the highest originating system ID will be preferred. In the case of overlapping keys in the winning South TIE, the behavior is undefined.

Consider that if a node goes down, nodes south of it will lose associated adjacencies causing them to disregard corresponding KVs. New KV South TIEs are advertised to prevent stale information being used by nodes that are farther south. KV advertisements southbound are not a result of independent computation by every node over the same set of South TIEs, but a diffused computation.

#### 4.3.4.2. Northbound

Certain use cases necessitate distribution of essential KV information that is generated by the leaves in the northbound direction. Such information is flooded in KV North TIEs. Since the originator of the KV North TIEs is preserved during flooding, overlapping keys MAY be used. However, to avoid further protocol complexity, the same tie-breaking rules as used in southbound distribution SHOULD be used.

#### 4.3.5. Interactions with BFD

RIFT MAY incorporate BFD [RFC5881] to react quickly to link failures. In such case following procedures are introduced:

After RIFT three-way hello adjacency convergence a BFD session MAY be formed automatically between the RIFT endpoints without further configuration using the exchanged discriminators. The capability of the remote side to support BFD is carried in the LIEs.

In case established BFD session goes Down after it was Up, RIFT adjacency SHOULD be re-initialized and subsequently started from Init after it sees a consecutive BFD Up.

In case of parallel links between nodes each link MAY run its own independent BFD session or they MAY share a session.

If link identifiers or BFD capabilities change, both the LIE and any BFD sessions SHOULD be brought down and back up again. In

case only the advertised capabilities change, the node MAY choose to persist the BFD session.

Multiple RIFT instances MAY choose to share a single BFD session, in such cases the behavior for which discriminators are used is undefined. However, RIFT MAY advertise the same link ID for the same interface in multiple instances to "share" discriminators.

BFD TTL follows [RFC5082].

#### 4.3.6. Fabric Bandwidth Balancing

A well understood problem in fabrics is that in case of link failures, it would be ideal to rebalance how much traffic is sent to switches in the next level based on available ingress and egress bandwidth.

RIFT supports a very light weight mechanism that can deal with the problem in an approximate way based on the fact that RIFT is loop-free.

##### 4.3.6.1. Northbound Direction

Every RIFT node SHOULD compute the amount of northbound bandwidth available through neighbors at higher level and modify distance received on default route from this neighbor. Default routes with differing distances SHOULD be used to support weighted ECMP forwarding. We call such a distance Bandwidth Adjusted Distance or BAD. This is best illustrated by a simple example.

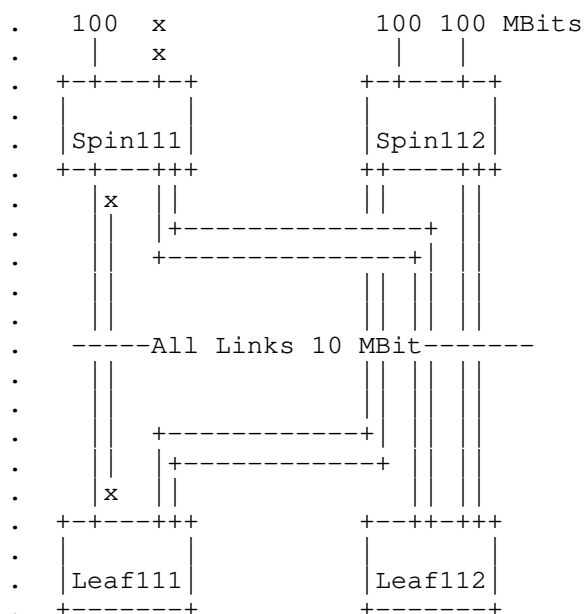


Figure 29: Balancing Bandwidth

Figure 29 depicts an example topology where links between leaf and spine nodes are 10 MBit/s and links from spine nodes northbound are 100 MBit/s. Consider a parallel link failure between Leaf 111 and Spine 111 and as a result, Leaf 111 wants to forward more traffic toward Spine 112. Additionally, we consider an uplink failure on Spine 111.

The local modification of the received default route distance from upper level is achieved by running a relatively simple algorithm where the bandwidth is weighted exponentially, while the distance on the default route represents a multiplier for the bandwidth weight for easy operational adjustments.

On a node, L, use Node TIEs to compute from each non-overloaded northbound neighbor N to compute 3 values:

L\_N\_u: as sum of the bandwidth available to N

N\_u: as sum of the uplink bandwidth available on N

T\_N\_u: as sum of  $L\_N\_u * OVERSUBSCRIPTION\_CONSTANT + N\_u$

For all  $T_{N_u}$  determine the according  $M_{N_u}$  as  $\log_2(\text{next\_power\_2}(T_{N_u}))$  and determine  $\text{MAX}_{M_{N_u}}$  as maximum value of all such  $M_{N_u}$  values.

For each advertised default route from a node N modify the advertised distance D to  $\text{BAD} = D * (1 + \text{MAX}_{M_{N_u}} - M_{N_u})$  and use BAD instead of distance D to weight balance default forwarding towards N.

For the example above, a simple table of values will help in understanding of the concept. We assume that all default route distances are advertised with  $D=1$  and that  $\text{OVERSUBSCRIPTION\_CONSTANT} = 1$ .

Node	N	$T_{N_u}$	$M_{N_u}$	BAD
Leaf111	Spine 111	110	7	2
Leaf111	Spine 112	220	8	1
Leaf112	Spine 111	120	7	2
Leaf112	Spine 112	220	8	1

Table 5: BAD Computation

If a calculation produces a result exceeding the range of the type, e.g. bandwidth, the result is set to the highest possible value for that type.

BAD SHOULD be only computed for default routes. A node MAY compute and use BAD for any disaggregated prefixes or other RIFT routes. A node MAY use a different algorithm to weight northbound traffic based on bandwidth. If a different algorithm is used, its successful behavior MUST NOT depend on uniformity of algorithm or synchronization of BAD computations across the fabric. E.g. it is conceivable that leaves could use real time link loads gathered by analytics to change the amount of traffic assigned to each default route next hop.

Furthermore, a change in available bandwidth will only affect, at most, two levels down in the fabric, i.e. the blast radius of bandwidth adjustments is constrained no matter the fabric's height.

#### 4.3.6.2. Southbound Direction

Due to its loop free nature, during South SPF, a node MAY account for maximum available bandwidth on nodes in lower levels and modify the amount of traffic offered to the next level's southbound nodes. It is worth considering that such computations may be more effective if standardized, but do not have to be. As long as a packet continues to flow southbound, it will take some viable, loop-free path to reach its destination.

#### 4.3.7. Label Binding

A node MAY advertise in its LIEs, a locally significant, downstream assigned, interface specific label. One use of such a label is a hop-by-hop encapsulation allowing forwarding planes to be easily distinguished among multiple RIFT instances.

#### 4.3.8. Leaf to Leaf Procedures

RIFT implementations SHOULD support special East-West adjacencies between leaf nodes. Leaf nodes supporting these procedures MUST:

- advertise the LEAF\_2\_LEAF flag in its node capabilities AND

- set the overload bit on all leaf's node TIEs AND

- flood only a node's own north and south TIEs over E-W leaf adjacencies AND

- always use E-W leaf adjacency in all SPF computations AND

- install a discard route for any advertised aggregate routes in a leaf's TIE AND

- never form southbound adjacencies.

This will allow the E-W leaf nodes to exchange traffic strictly for the prefixes advertised in each other's north prefix TIEs (since the southbound computation will find the reverse direction in the other node's TIE and install its north prefixes).

#### 4.3.9. Address Family and Multi Topology Considerations

Multi-Topology (MT) [RFC5120] and Multi-Instance (MI) [RFC8202] concepts are used today in link-state routing protocols to support several domains on the same physical topology. RIFT supports this capability by carrying transport ports in the LIE protocol exchanges.

Multiplexing of LIEs can be achieved by either choosing varying multicast addresses or ports on the same address.

BFD interactions in Section 4.3.5 are implementation dependent when multiple RIFT instances run on the same link.

#### 4.3.10. Reachability of Internal Nodes in the Fabric

RIFT does not require that nodes have reachable addresses in the fabric, though it is clearly desirable for operational purposes. Under normal operating conditions this can be easily achieved by injecting the node's loopback address into North and South Prefix TIEs or other implementation specific mechanisms.

Special considerations arise when a node loses all northbound adjacencies, but is not at the top of the fabric. These are outside the scope of this document and could be discussed in a separate document.

#### 4.3.11. One-Hop Healing of Levels with East-West Links

Based on the rules defined in Section 4.2.4, Section 4.2.3.8 and given presence of E-W links, RIFT can provide a one-hop protection for nodes that lost all their northbound links. This can also be applied to multi-plane designs where complex link set failures occur at the Top-of-Fabric when links are exclusively used for flooding topology information. Section 5.4 outlines this behavior.

### 4.4. Security

#### 4.4.1. Security Model

An inherent property of any security and ZTP architecture is the resulting trade-off in regard to integrity verification of the information distributed through the fabric vs. provisioning and auto-configuration requirements. At a minimum the security of an established adjacency should be ensured. The stricter the security model the more provisioning must take over the role of ZTP.

RIFT supports the following security models to allow for flexible control by the operator.

- o The most security conscious operators may choose to have control over which ports interconnect between a given pair of nodes, we call this the "Port-Association Model" (PAM). This is achievable by configuring each pair of directly connected ports with a designated shared key or public/private key pair.



- o In physically secure data center locations, operators may choose to control connectivity between entire nodes, we call this the "Node-Association Model" (NAM). A benefit of this model is that it allows for simplified port sparing.
- o In the most relaxed environments, an operator may only choose to control which nodes join a particular fabric. We call this the "Fabric-Association Model" (FAM). This is achievable by using a single shared secret across the entire fabric. Such flexibility makes sense when we consider servers as leaf devices, which are replaced more often than network nodes. In addition, this model allows for simplified node sparing.
- o These models may be mixed throughout the fabric depending upon security requirements at various levels of the fabric and willingness to accept increased provisioning complexity.

In order to support the cases mentioned above, RIFT implementations supports, through operator control, mechanisms that allow for:

- a. specification of the appropriate level in the fabric,
- b. discovery and reporting of missing connections,
- c. discovery and reporting of unexpected connections while preventing them from forming insecure adjacencies.

Operators may only choose to configure the level of each node, but not explicitly configure which connections are allowed. In this case, RIFT will only allow adjacencies to establish between nodes that are in adjacent levels. Operators with the lowest security requirements may not use any configuration to specify which connections are allowed. Nodes in such fabrics could rely fully on ZTP and only established adjacencies between nodes in adjacent levels. Figure 30 illustrates inherent tradeoffs between the different security models.

Some level of link quality verification may be required prior to an adjacency being used for forwarding. For example, an implementation may require that a BFD session comes up before advertising the adjacency.

For the cases outlined above, RIFT has two approaches to enforce that a local port is connected to the correct port on the correct remote node. One approach is to piggy-back on RIFT's authentication mechanism. Assuming the provisioning model (e.g. the YANG model) is flexible enough, operators can choose to provision a unique authentication key for:

- a. each pair of ports in "port-association model" or
- b. each pair of switches in "node-association model" or
- c. each pair of levels or
- d. the entire fabric in "fabric-association model".

The other approach is to rely on the system-id, port-id and level fields in the LIE message to validate an adjacency against the expected cabling topology, and optionally introduce some new rules in the FSM to allow the adjacency to come up if the expectations are met.

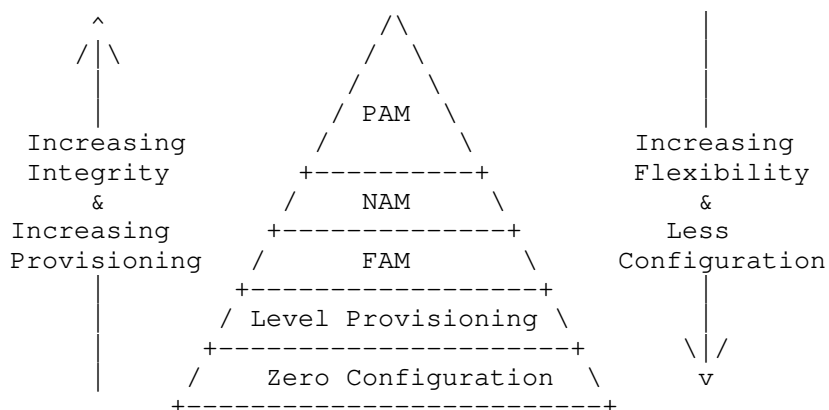


Figure 30: Security Model

#### 4.4.2. Security Mechanisms

RIFT Security goals are to ensure:

1. authentication
2. message integrity
3. the prevention of replay attacks
4. low processing overhead
5. efficient messaging

Message confidentiality is a non-goal.

The model in the previous section allows a range of security key types that are analogous to the various security association models. PAM and NAM allow security associations at the port or node level using symmetric or asymmetric keys that are pre-installed. FAM argues for security associations to be applied only at a group level or to be refined once the topology has been established. RIFT does not specify how security keys are installed or updated, though it does specify how the key can be used to achieve security goals.

The protocol has provisions for "weak" nonces to prevent replay attacks and includes authentication mechanisms comparable to [RFC5709] and [RFC7987].

#### 4.4.3. Security Envelope

RIFT MUST be carried in a mandatory secure envelope illustrated in Figure 31. Any value in the packet following a security fingerprint MUST be used only after the appropriate fingerprint has been validated.

Local configuration MAY allow for the envelope's integrity checks to be skipped.

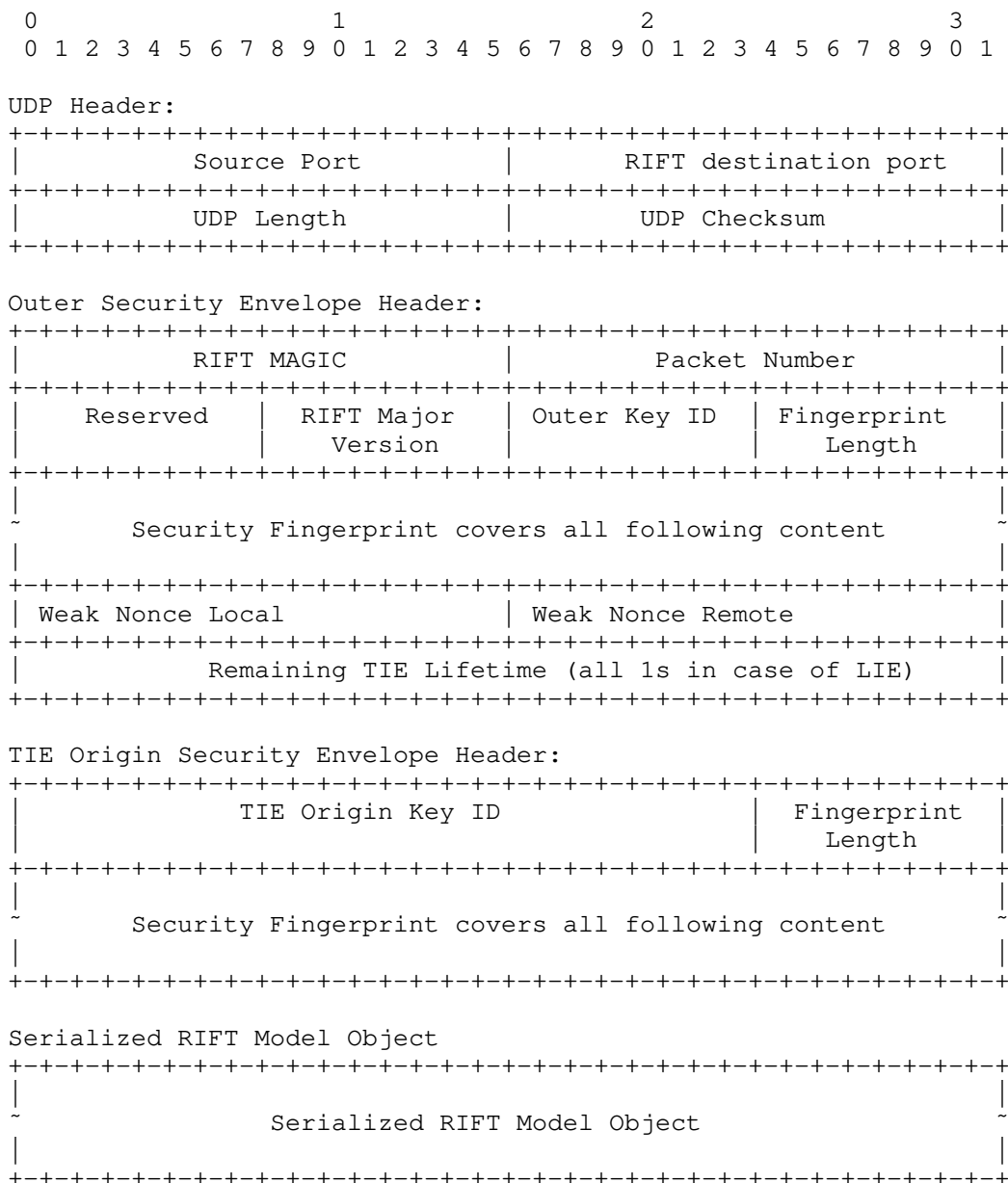


Figure 31: Security Envelope

RIFT MAGIC: 16 bits. Constant value of 0xA1F7 that allows to classify RIFT packets independent of used UDP port.

Packet Number: 16 bits. An optional, per packet type monotonically growing number rolling over using sequence number arithmetic defined in Appendix A. A node SHOULD correctly set the number on subsequent packets or otherwise MUST set the value to `'undefined_packet_number'` as provided in the schema. This number can be used to detect losses and misordering in flooding for either operational purposes or in implementation to adjust flooding behavior to current link or buffer quality. This number MUST NOT be used to discard or validate the correctness of packets.

RIFT Major Version: 8 bits. It allows to check whether protocol versions are compatible, i.e. if the serialized object can be decoded at all. An implementation MUST drop packets with unexpected values and MAY report a problem.

Outer Key ID: 8 bits to allow key rollovers. This implies key type and algorithm. Value 0 means that no valid fingerprint was computed. This key ID scope is local to the nodes on both ends of the adjacency.

TIE Origin Key ID: 24 bits. This implies key type and used algorithm. Value 0 means that no valid fingerprint was computed. This key ID scope is global to the RIFT instance since it implies the originator of the TIE so the contained object does not have to be de-serialized to obtain it.

Length of Fingerprint: 8 bits. Length in 32-bit multiples of the following fingerprint (not including lifetime or weak nonces). It allows the structure to be navigated when an unknown key type is present. To clarify, a common corner case when this value is set to 0 is when it signifies an empty (0 bytes long) security fingerprint.

Security Fingerprint: 32 bits \* Length of Fingerprint. This is a signature that is computed over all data following after it. If the significant bits of fingerprint are fewer than the 32 bits padded length than the significant bits MUST be left aligned and remaining bits on the right padded with 0s. When using PKI the Security fingerprint originating node uses its private key to create the signature. The original packet can then be verified provided the public key is shared and current.

Remaining TIE Lifetime: 32 bits. In case of anything but TIEs this field MUST be set to all ones and Origin Security Envelope Header

MUST NOT be present in the packet. For TIEs this field represents the remaining lifetime of the TIE and Origin Security Envelope Header MUST be present in the packet. The value in the serialized model object MUST be ignored.

Weak Nonce Local: 16 bits. Local Weak Nonce of the adjacency as advertised in LIEs.

Weak Nonce Remote: 16 bits. Remote Weak Nonce of the adjacency as received in LIEs.

TIE Origin Security Envelope Header: It MUST be present if and only if the Remaining TIE Lifetime field is NOT all ones. It carries through the originators key ID and according fingerprint of the object to protect TIE from modification during flooding. This ensures origin validation and integrity (but does not provide validation of a chain of trust).

Observe that due to the schema migration rules per Appendix B the contained model can be always decoded if the major version matches and the envelope integrity has been validated. Consequently, description of the TIE is available to flood it properly including unknown TIE types.

#### 4.4.4. Weak Nonces

The protocol uses two 16 bit nonces to salt generated signatures. We use the term "nonce" a bit loosely since RIFT nonces are not being changed in every packet as common in cryptography. For efficiency purposes they are changed at a high enough frequency to dwarf practical replay attack attempts. Therefore, we call them "weak" nonces.

Any implementation including RIFT security MUST generate and wrap around local nonces properly. When a nonce increment leads to 'undefined\_nonce' value, the value MUST be incremented again immediately. All implementation MUST reflect the neighbor's nonces. An implementation SHOULD increment a chosen nonce on every LIE FSM transition that ends up in a different state from the previous and MUST increment its nonce at least every 5 minutes (such considerations allow for efficient implementations without opening a significant security risk). When flooding TIEs, the implementation MUST use recent (i.e. within allowed difference) nonces reflected in the LIE exchange. The schema specifies the maximum allowable nonce value difference on a packet compared to reflected nonces in the LIEs. Any packet received with nonces deviating more than the allowed delta MUST be discarded without further computation of signatures to prevent computation load attacks.

In cases where a secure implementation does not receive signatures or receives undefined nonces from a neighbor (indicating that it does not support or verify signatures), it is a matter of local policy as to how those packets are treated. A secure implementation MAY refuse forming an adjacency with an implementation that is not advertising signatures or valid nonces, or it MAY continue signing local packets while accepting a neighbor's packets without further security validation.

As a necessary exception, an implementation MUST advertise the remote nonce value as 'undefined\_nonce' when the FSM is not in two-way or three-way state and accept an 'undefined\_nonce' for its local nonce value on packets in any other state than three-way.

As optional optimization, an implementation MAY send one LIE with previously negotiated neighbor's nonce to try to speed up a neighbor's transition from three-way to one-way and MUST revert to sending 'undefined\_nonce' after that.

#### 4.4.5. Lifetime

Protecting flooding lifetime may lead to an excessive number of security fingerprint computations and to avoid this the application generating the fingerprints for advertised TIEs, MAY round the value down to the next 'rounddown\_lifetime\_interval'. Such an optimization in the presence of security hashes over advancing weak nonces, may not be feasible.

#### 4.4.6. Key Management

As outlined in Section Section 7, either a private shared key or a public/private key pair is used to authenticate the adjacency. Both the key distribution and key synchronization methods are out of scope for this document. Both nodes in the adjacency MUST share the same keys, key type, and algorithm for a given key ID. Mismatched keys will not inter-operate as their security envelopes will be unverifiable.

Key roll-over while the adjacency is active MAY be supported. The specific mechanism is well documented in [RFC6518].

#### 4.4.7. Security Association Changes

There is no mechanism to convert a security envelope for the same key ID from one algorithm to another once the envelope is operational. The recommended procedure to change to a new algorithm is to take the adjacency down, make the necessary changes, and bring the adjacency back up. Obviously, an implementation MAY choose to stop verifying

security envelope for the duration of algorithm change to keep the adjacency up but since this introduces a security vulnerability window, such roll-over SHOULD NOT be recommended.

5. Examples

5.1. Normal Operation

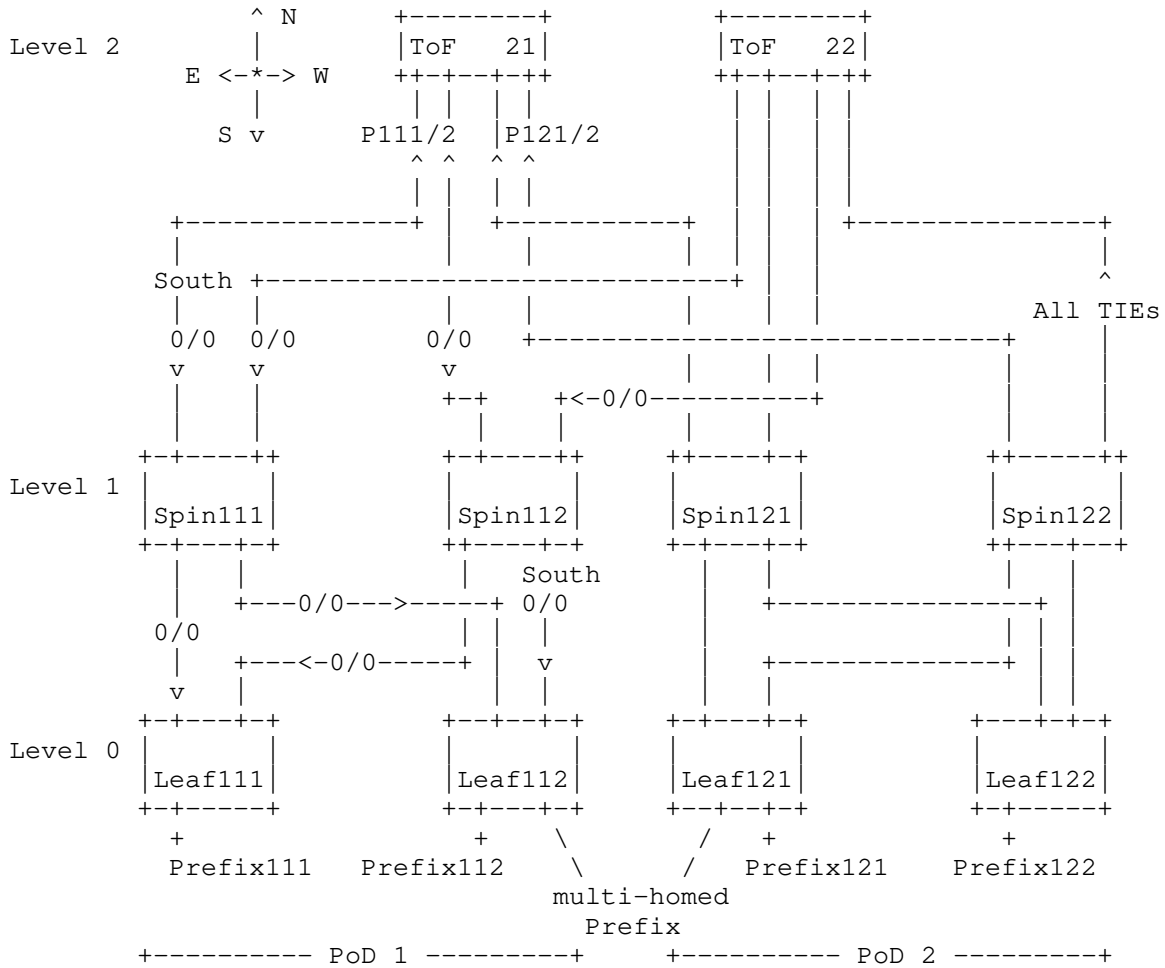


Figure 32: Normal Case Topology

This section describes RIFT deployment in example topology given in Figure 32 without any node or link failures. We disregard flooding



reduction for simplicity's sake and compress the node names in some cases to fit them into the picture better.

First, the following bi-directional adjacencies will be established:

1. ToF 21 (PoD 0) to Spine 111, Spine 112, Spine 121, and Spine 122
2. ToF 22 (PoD 0) to Spine 111, Spine 112, Spine 121, and Spine 122
3. Spine 111 to Leaf 111, Leaf 112
4. Spine 112 to Leaf 111, Leaf 112
5. Spine 121 to Leaf 121, Leaf 122
6. Spine 122 to Leaf 121, Leaf 122

Leaf 111 and Leaf 112 originate N-TIEs for Prefix 111 and Prefix 112 (respectively) to both Spine 111 and Spine 112 (Leaf 112 also originates an N-TIE for the multi-homed prefix). Spine 111 and Spine 112 will then originate their own N-TIEs, as well as flood the N-TIEs received from Leaf 111 and Leaf 112 to both ToF 21 and ToF 22.

Similarly, Leaf 121 and Leaf 122 originate North TIEs for Prefix 121 and Prefix 122 (respectively) to Spine 121 and Spine 122 (Leaf 121 also originates an North TIE for the multi-homed prefix). Spine 121 and Spine 122 will then originate their own North TIEs, as well as flood the North TIEs received from Leaf 121 and Leaf 122 to both ToF 21 and ToF 22.

Spines hold only North TIEs of level 0 for their PoD, while leaves only hold their own North TIEs while at this point, both ToF 21 and ToF 22 (as well as any northbound connected controllers) would have the complete network topology.

ToF 21 and ToF 22 would then originate and flood South TIEs containing any established adjacencies and a default IP route to all spines. Spine 111, Spine 112, Spine 121, and Spine 122 will reflect all Node South TIEs received from ToF 21 to ToF 22, and all Node South TIEs from ToF 22 to ToF 21. South TIEs will not be re-propagated southbound.

South TIEs containing a default IP route are then originated by both Spine 111 and Spine 112 toward Leaf 111 and Leaf 112. Similarly, South TIEs containing a default IP route are originated by Spine 121 and Spine 122 toward Leaf 121 and Leaf 122.

At this point IP connectivity across maximum number of viable paths has been established for all leaves, with routing information constrained to only the minimum amount that allows for normal operation and redundancy.

5.2. Leaf Link Failure

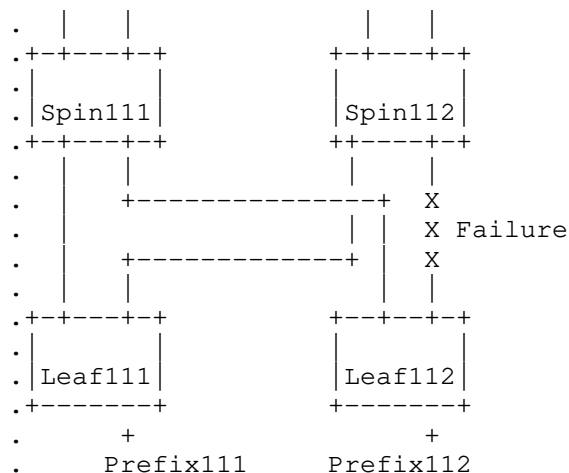


Figure 33: Single Leaf Link Failure

In the event of a link failure between Spine 112 and Leaf 112, both nodes will originate new Node TIEs that contain their connected adjacencies, except for the one that just failed. Leaf 112 will send a Node North TIE to Spine 111. Spine 112 will send a Node North TIE to ToF 21 and ToF 22 as well as a new Node South TIE to Leaf 111 that will be reflected to Spine 111. Necessary SPF recomputation will occur, resulting in Spine 112 no longer being in the forwarding path for Prefix 112.

Spine 111 will also disaggregate Prefix 112 by sending new Prefix South TIE to Leaf 111 and Leaf 112. Though we cover disaggregation in more detail in the following section, it is worth mentioning in this example as it further illustrates RIFT's blackhole mitigation mechanism. Consider that Leaf 111 has yet to receive the more specific (disaggregated) route from Spine 111. In such a scenario, traffic from Leaf 111 toward Prefix 112 may still use Spine 112's default route, causing it to traverse ToF 21 and ToF 22 back down via Spine 111. While this behavior is suboptimal, it is transient in nature and preferred to black-holing traffic.

5.3. Partitioned Fabric

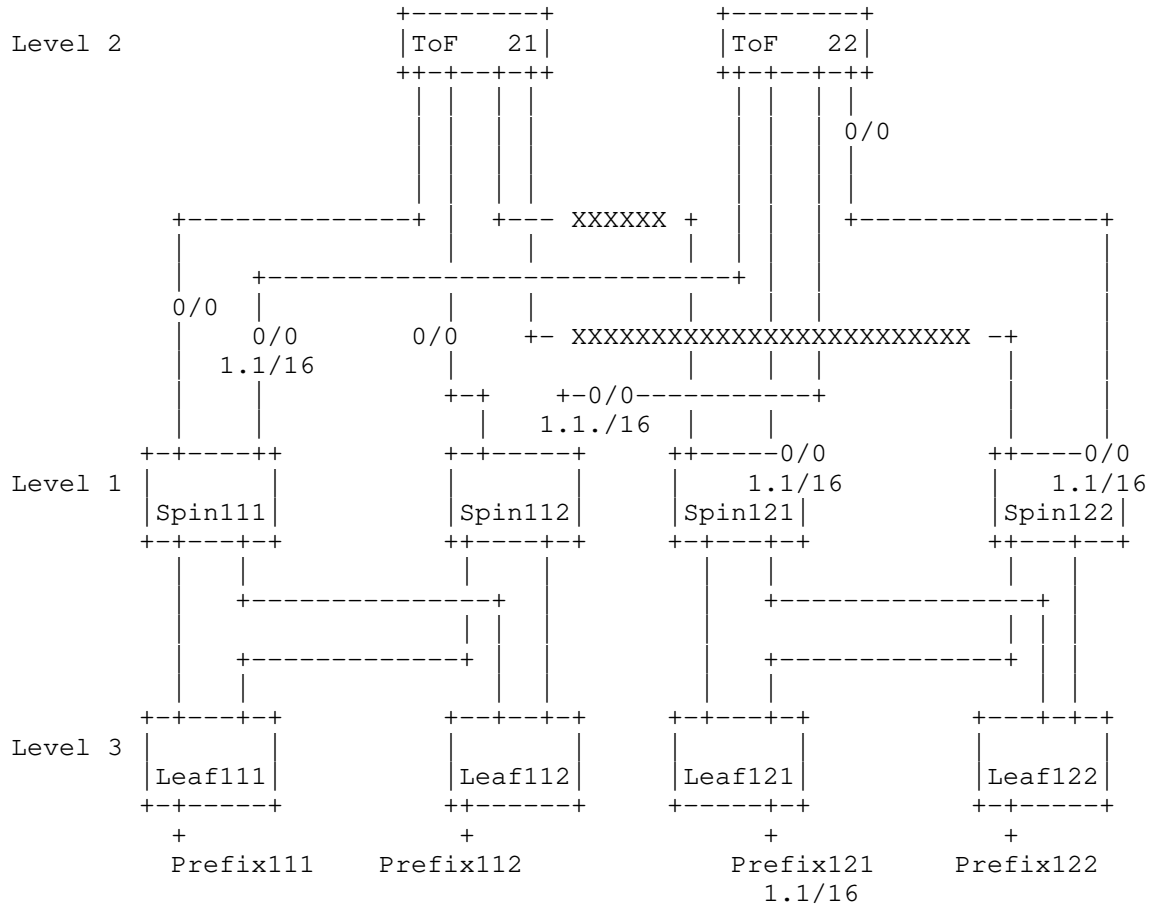


Figure 34: Fabric Partition

Figure 34 shows one of more catastrophic scenarios where ToF 21 is completely severed from access to Prefix 121 due to a double link failure. If only default routes existed, this would result in 50% of traffic from Leaf 111 and Leaf 112 toward Prefix 121 being black-holed.

The mechanism to resolve this scenario hinges on ToF 21's South TIEs being reflected from Spine 111 and Spine 112 to ToF 22. Once ToF 22 sees that Prefix 121 cannot be reached from ToF 21, it will begin to disaggregate Prefix 121 by advertising a more specific route (1.1/16)

along with the default IP prefix route to all spines (ToF 21 still only sends a default route). The result is Spine 111 and Spine112 using the more specific route to Prefix 121 via ToF 22. All other prefixes continue to use the default IP prefix route toward both ToF 21 and ToF 22.

The more specific route for Prefix 121 being advertised by ToF 22 does not need to be propagated further south to the leaves, as they do not benefit from this information. Spine 111 and Spine 112 are only required to reflect the new South Node TIEs received from ToF 22 to ToF 21. In short, only the relevant nodes received the relevant updates, thereby restricting the failure to only the partitioned level rather than burdening the whole fabric with the flooding and recomputation of the new topology information.

To finish our example, the following table shows sets computed by ToF 22 using notation introduced in Section 4.2.5:

|R = Prefix 111, Prefix 112, Prefix 121, Prefix 122

|H (for r=Prefix 111) = Spine 111, Spine 112

|H (for r=Prefix 112) = Spine 111, Spine 112

|H (for r=Prefix 121) = Spine 121, Spine 122

|H (for r=Prefix 122) = Spine 121, Spine 122

|A (for ToF 21) = Spine 111, Spine 112

With that and |H (for r=Prefix 121) and |H (for r=Prefix 122) being disjoint from |A (for ToF 21), ToF 22 will originate an South TIE with Prefix 121 and Prefix 122, which will be flooded to all spines.

#### 5.4. Northbound Partitioned Router and Optional East-West Links

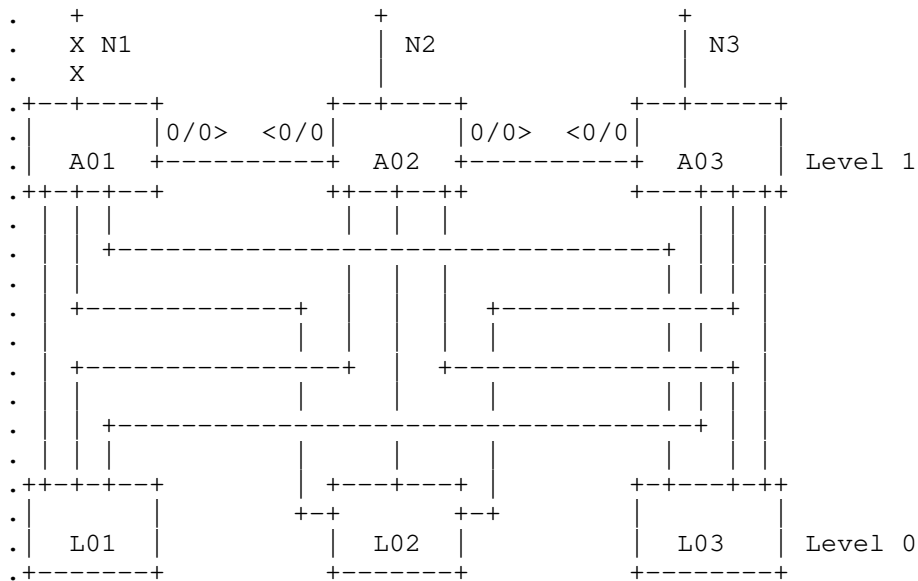


Figure 35: North Partitioned Router

Figure 35 shows a part of a fabric where level 1 is horizontally connected and A01 lost its only northbound adjacency. Based on N-SPF rules in Section 4.2.4.1 A01 will compute northbound reachability by using the link A01 to A02. A02 however, will NOT use this link during N-SPF. The result is A01 utilizing the horizontal link for default route advertisement and unidirectional routing.

Furthermore, if A02 also loses its only northbound adjacency (N2), the situation evolves. A01 will no longer have northbound reachability while it sees A03's northbound adjacencies in South Node TIEs reflected by nodes south of it. As a result, A01 will no longer advertise its default route in accordance with Section 4.2.3.8.

## 6. Implementation and Operation: Further Details

### 6.1. Considerations for Leaf-Only Implementation

RIFT can and is intended to be stretched to the lowest level in the IP fabric to integrate ToRs or even servers. Since those entities would run as leaves only, it is worth to observe that a leaf only version is significantly simpler to implement and requires much less resources:

1. Leaf nodes only need to maintain a multipath default route under normal circumstances. However, in cases of catastrophic partitioning, leaf nodes SHOULD be capable of accommodating all the leaf routes in its own PoD to prevent black-holing.
2. Leaf nodes hold only their own North TIEs and South TIEs of Level 1 nodes they are connected to.
3. Leaf nodes do not have to support any type of de-aggregation computation or propagation.
4. Leaf nodes are not required to support overload bit.
5. Leaf nodes do not need to originate S-TIEs unless optional leaf-2-leaf features are desired.

### 6.2. Considerations for Spine Implementation

Spine nodes will never act as Top of Fabric, and are therefore not required to run a full RIFT implementation. Specifically, spines do not need to perform negative disaggregation computation other than respecting northbound disaggregation advertised from the north.

### 6.3. Adaptations to Other Proposed Data Center Topologies

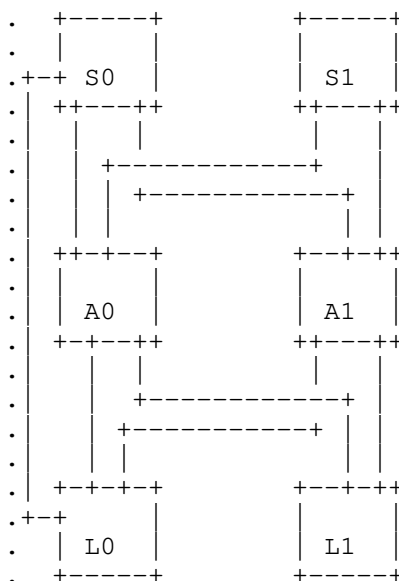


Figure 36: Level Shortcut

RIFT is not strictly limited to Clos topologies. The protocol only requires a sense of "compass rose directionality" either achieved through configuration or derivation of levels. So, conceptually, leaf-2-leaf links and even shortcuts between levels could be included. Figure 36 depicts an example of a shortcut between levels. In this example, sub-optimal routing will occur when traffic is sent from L0 to L1 via S0's default route and back down through A0 or A1. In order to ensure that only default routes from A0 or A1 are used, all leaves would be required to install each others routes.

While various technical and operational challenges may require the use of such modifications, discussion of those topics are outside the scope of this document.

#### 6.4. Originating Non-Default Route Southbound

An implementation MAY choose to originate more specific prefixes (P') southbound instead of only the default route (as described in Section 4.2.3.8). In such a scenario, all addresses carried within the RIFT domain MUST be contained within P'.

### 7. Security Considerations

#### 7.1. General

One can consider attack vectors where a router may reboot many times while changing its system ID and pollute the network with many stale TIEs or TIEs are sent with very long lifetimes and not cleaned up when the routes vanish. Those attack vectors are not unique to RIFT. Given large memory footprints available today those attacks should be relatively benign. Otherwise a node SHOULD implement a strategy of discarding contents of all TIEs that were not present in the SPF tree over a certain, configurable period of time. Since the protocol, like all modern link-state protocols, is self-stabilizing and will advertise the presence of such TIEs to its neighbors, they can be re-requested again if a computation finds that it sees an adjacency formed towards the system ID of the discarded TIEs.

#### 7.2. ZTP

Section 4.2.7 presents many attack vectors in untrusted environments, starting with nodes that oscillate their level offers to the possibility of nodes offering a three-way adjacency with the highest possible level value and a very long holdtime trying to put itself "on top of the lattice" thereby allowing it to gain access to the whole southbound topology. Session authentication mechanisms are necessary in environments where this is possible and RIFT provides the security envelope to ensure this if so desired.

### 7.3. Lifetime

Traditional IGP protocols are vulnerable to lifetime modification and replay attacks that can be somewhat mitigated by using techniques like [RFC7987]. RIFT removes this attack vector by protecting the lifetime behind a signature computed over it and additional nonce combination which makes even the replay attack window very small and for practical purposes irrelevant since lifetime cannot be artificially shortened by the attacker.

### 7.4. Packet Number

Optional packet number is carried in the security envelope without any encryption protection and is hence vulnerable to replay and modification attacks. Contrary to nonces this number must change on every packet and would present a very high cryptographic load if signed. The attack vector packet number present is relatively benign. Changing the packet number by a man-in-the-middle attack will only affect operational validation tools and possibly some performance optimizations on flooding. It is expected that an implementation detecting too many "fake losses" or "misorderings" due to the attack on the packet number would simply suppress its further processing.

### 7.5. Outer Fingerprint Attacks

A node can try to inject LIE packets observing a conversation on the wire by using the outer key ID albeit it cannot generate valid hashes in case it changes the integrity of the message so the only possible attack is DoS due to excessive LIE validation.

A node can try to replay previous LIEs with changed state that it recorded but the attack is hard to replicate since the nonce combination must match the ongoing exchange and is then limited to a single flap only since both nodes will advance their nonces in case the adjacency state changed. Even in the most unlikely case the attack length is limited due to both sides periodically increasing their nonces.

### 7.6. TIE Origin Fingerprint DoS Attacks

A compromised node can attempt to generate "fake TIEs" using other nodes' TIE origin key identifiers. Albeit the ultimate validation of the origin fingerprint will fail in such scenarios and not progress further than immediately peering nodes, the resulting denial of service attack seems unavoidable since the TIE origin key id is only protected by the, here assumed to be compromised, node.



### 7.7. Host Implementations

It can be reasonably expected that with the proliferation of RotH servers, rather than dedicated networking devices, will represent a significant amount of RIFT devices. Given their normally far wider software envelope and access granted to them, such servers are also far more likely to be compromised and present an attack vector on the protocol. Hijacking of prefixes to attract traffic is a trust problem and cannot be easily addressed within the protocol if the trust model is breached, i.e. the server presents valid credentials to form an adjacency and issue TIEs. In an even more devious way, the servers can present DoS (or even DDos) vectors of issuing too many LIE packets, flood large amounts of North TIEs and attempt similar resource overrun attacks. A prudent implementation forming adjacencies to leaves should implement according thresholds mechanisms and raise warnings when e.g. a leaf is advertising an excess number of TIEs or prefixes. Additionally, such implementation could refuse any topology information except the node's own TIEs and authenticated, reflected South Node TIEs at own level.

To isolate possible attack vectors on the leaf to the largest possible extent a dedicated leaf-only implementation could run without any configuration by hard-coding a well-known adjacency key (which can be always rolled-over by the means of e.g. well-known key-value distributed from top of the fabric), leaf level value and always setting overload bit. All other values can be derived by automatic means as described earlier in the protocol specification.

## 8. IANA Considerations

This specification requests multicast address assignments and standard port numbers. Additionally registries for the schema are requested and suggested values provided that reflect the numbers allocated in the given schema.

### 8.1. Requested Multicast and Port Numbers

This document requests allocation in the 'IPv4 Multicast Address Space' registry the suggested value of 224.0.0.120 as 'ALL\_V4\_RIFT\_ROUTERS' and in the 'IPv6 Multicast Address Space' registry the suggested value of FF02::A1F7 as 'ALL\_V6\_RIFT\_ROUTERS'.

This document requests allocation in the 'Service Name and Transport Protocol Port Number Registry' the allocation of a suggested value of 914 on udp for 'RIFT\_LIES\_PORT' and suggested value of 915 for 'RIFT\_TIES\_PORT'.

## 8.2. Requested Registries with Suggested Values

This section requests registries that help govern the schema via usual IANA registry procedures. A top level 'RIFT' registry should hold the according registries requested in following sections with their pre-defined values. IANA is requested to store the schema version introducing the allocated value as well as, optionally, its description when present. This will allow to assign different values to an entry depending on schema version. Alternately, IANA is requested to consider a root RIFT/3 registry to store RIFT schema major version 3 values and may be requested in the future to create a RIFT/4 registry under that. In any case, IANA is requested to store the schema version in the entries since that will allow to distinguish between minor versions in the same major schema version. All values not suggested as to be considered 'Unassigned'. The range of every registry is a 16-bit integer. Allocation of new values is always performed via 'Expert Review' action.

### 8.2.1. Registry RIFT\_v4/common/AddressFamilyType

Address family type.

#### 8.2.1.1. Requested Entries

Name	Value	Schema	Version	Description
Illegal	0		4.1	
AddressFamilyMinValue	1		4.1	
IPv4	2		4.1	
IPv6	3		4.1	
AddressFamilyMaxValue	4		4.1	

### 8.2.2. Registry RIFT\_v4/common/HierarchyIndications

Flags indicating node configuration in case of ZTP.

#### 8.2.2.1. Requested Entries

Name	Value	Schema	Version	Description
leaf_only	0		4.1	
leaf_only_and_leaf_2_leaf_procedures	1		4.1	
top_of_fabric	2		4.1	

### 8.2.3. Registry RIFT\_v4/common/IEEE802\_1ASTimeStampType

Timestamp per IEEE 802.1AS, all values MUST be interpreted in implementation as unsigned.

## 8.2.3.1. Requested Entries

Name	Value	Schema	Version	Description
AS_sec	1		4.1	
AS_nsec	2		4.1	

## 8.2.4. Registry RIFT\_v4/common/IPAddressType

IP address type.

## 8.2.4.1. Requested Entries

Name	Value	Schema	Version	Description
ipv4address	1		4.1	Content is IPv4
ipv6address	2		4.1	Content is IPv6

## 8.2.5. Registry RIFT\_v4/common/IPPrefixType

Prefix advertisement.

@note: for interface addresses the protocol can propagate the address part beyond the subnet mask and on reachability computation that has to be normalized. The non-significant bits can be used for operational purposes.

## 8.2.5.1. Requested Entries

Name	Value	Schema	Version	Description
ipv4prefix	1		4.1	
ipv6prefix	2		4.1	

## 8.2.6. Registry RIFT\_v4/common/IPv4PrefixType

IPv4 prefix type.

## 8.2.6.1. Requested Entries

Name	Value	Schema	Version	Description
address	1		4.1	
prefixlen	2		4.1	

## 8.2.7. Registry RIFT\_v4/common/IPv6PrefixType

IPv6 prefix type.

## 8.2.7.1. Requested Entries

Name	Value	Schema	Version	Description
address	1		4.1	
prefixlen	2		4.1	

## 8.2.8. Registry RIFT\_v4/common/PrefixSequenceType

Sequence of a prefix in case of move.

## 8.2.8.1. Requested Entries

Name	Value	Schema	Description
		Version	
timestamp	1	4.1	
transactionid	2	4.1	Transaction ID set by client in e.g. in 6LoWPAN.

## 8.2.9. Registry RIFT\_v4/common/RouteType

RIFT route types.

@note: route types which MUST be ordered on their preference PGP prefixes are most preferred attracting traffic north (towards spine) and then south normal prefixes are attracting traffic south (towards leafs), i.e. prefix in NORTH PREFIX TIE is preferred over SOUTH PREFIX TIE.

@note: The only purpose of those values is to introduce an ordering whereas an implementation can choose internally any other values as long the ordering is preserved

## 8.2.9.1. Requested Entries

Name	Value	Schema	Version	Description
Illegal	0		4.1	
RouteTypeMinValue	1		4.1	
Discard	2		4.1	
LocalPrefix	3		4.1	
SouthPGPPrefix	4		4.1	
NorthPGPPrefix	5		4.1	
NorthPrefix	6		4.1	
NorthExternalPrefix	7		4.1	
SouthPrefix	8		4.1	
SouthExternalPrefix	9		4.1	
NegativeSouthPrefix	10		4.1	
RouteTypeMaxValue	11		4.1	

## 8.2.10. Registry RIFT\_v4/common/TIETypeType

Type of TIE.

This enum indicates what TIE type the TIE is carrying. In case the value is not known to the receiver, the TIE MUST be re-flooded. This allows for future extensions of the protocol within the same major schema with types opaque to some nodes UNLESS the flooding scope is not the same as prefix TIE, then a major version revision MUST be performed.

## 8.2.10.1. Requested Entries

Name	Value	Schema Version	Description
Illegal	0	4.1	
TIETypeMinValue	1	4.1	
NodeTIEType	2	4.1	
PrefixTIEType	3	4.1	
PositiveDisaggregationPrefixTIEType	4	4.1	
NegativeDisaggregationPrefixTIEType	5	4.1	
PGPrefixTIEType	6	4.1	
KeyValueTIEType	7	4.1	
ExternalPrefixTIEType	8	4.1	
PositiveExternalDisaggregationPrefixTIEType	9	4.1	
TIETypeMaxValue	10	4.1	

## 8.2.11. Registry RIFT\_v4/common/TieDirectionType

Direction of TIEs.

## 8.2.11.1. Requested Entries

Name	Value	Schema Version	Description
Illegal	0	4.1	
South	1	4.1	
North	2	4.1	
DirectionMaxValue	3	4.1	

## 8.2.12. Registry RIFT\_v4/encoding/Community

Prefix community.

## 8.2.12.1. Requested Entries

Name	Value	Schema Version	Description
top	1	4.1	Higher order bits
bottom	2	4.1	Lower order bits

## 8.2.13. Registry RIFT\_v4/encoding/KeyValueTIEElement

Generic key value pairs.

## 8.2.13.1. Requested Entries

Name	Value	Schema	Version	Description
keyvalues	1		4.1	

## 8.2.14. Registry RIFT\_v4/encoding/LIEPacket

RIFT LIE Packet.

@note: this node's level is already included on the packet header

## 8.2.14.1. Requested Entries

Name	Value	Schema	Description
		Version	
name	1	4.1	Node or adjacency name.
local_id	2	4.1	Local link ID.
flood_port	3	4.1	UDP port to which we can receive flooded TIEs.
link_mtu_size	4	4.1	Layer 3 MTU, used to discover to mismatch.
link_bandwidth	5	4.1	Local link bandwidth on the interface.
neighbor	6	4.1	Reflects the neighbor once received to provide 3-way connectivity.
pod	7	4.1	Node's PoD.
node_capabilities	10	4.1	Node capabilities shown in LIE. The capabilities MUST match the capabilities shown in the Node TIEs, otherwise the behavior is unspecified. A node detecting the mismatch SHOULD generate according error.
link_capabilities	11	4.1	Capabilities of this link.
holdtime	12	4.1	Required holdtime of the adjacency, i.e. how much time MUST expire without LIE for the adjacency to drop.
label	13	4.1	Unsolicited, downstream assigned locally

Name	Value	Schema Version	Description
not_a_ztp_offer	21	4.1	Indicates that the level on the LIE MUST NOT be used to derive a ZTP level by the receiving node.
you_are_flood_repeater	22	4.1	Indicates to northbound neighbor that it should be reflooding this node's N-TIEs to achieve flood reduction and balancing for northbound flooding. To be ignored if received from a northbound adjacency.
you_are_sending_too_quickly	23	4.1	Can be optionally set to indicate to neighbor that packet losses are seen on reception based on packet numbers or the rate is too high. The receiver SHOULD temporarily slow down flooding rates.
instance_name	24	4.1	Instance name in case multiple RIFT instances running on same interface.

#### 8.2.15. Registry RIFT\_v4/encoding/LinkCapabilities

Link capabilities.

##### 8.2.15.1. Requested Entries

Name	Value	Schema Version	Description
bfd	1	4.1	Indicates that the link is supporting BFD.
v4_forwarding_capable	2	4.1	Indicates whether the interface will support v4 forwarding.

#### 8.2.16. Registry RIFT\_v4/encoding/LinkIDPair

LinkID pair describes one of parallel links between two nodes.

## 8.2.16.1. Requested Entries

Name	Value	Schema	Description
		Version	
local_id	1	4.1	Node-wide unique value for the local link.
remote_id	2	4.1	Received remote link ID for this link.
platform_interface_index	10	4.1	Describes the local interface index of the link.
platform_interface_name	11	4.1	Describes the local interface name.
trusted_outer_security_key	12	4.1	Indication whether the link is secured, i.e. protected by outer key, absence of this element means no indication, undefined outer key means not secured.
bfd_up	13	4.1	Indication whether the link is protected by established BFD session.
address_families	14	4.1	Optional indication which address families are up on the interface

## 8.2.17. Registry RIFT\_v4/encoding/Neighbor

Neighbor structure.

## 8.2.17.1. Requested Entries

Name	Value	Schema	Version	Description
originator	1		4.1	System ID of the originator.
remote_id	2		4.1	ID of remote side of the link.

## 8.2.18. Registry RIFT\_v4/encoding/NodeCapabilities

Capabilities the node supports.

@note: The schema may add to this field future capabilities to indicate whether it will support interpretation of future schema extensions on the same major revision. Such fields MUST be optional and have an implicit or explicit false default value. If a future capability changes route selection or generates blackholes if some nodes are not supporting it then a major version increment is unavoidable.



## 8.2.18.1. Requested Entries

Name	Value	Schema Version	Description
protocol_minor_version	1	4.1	Must advertise supported minor version dialect that way.
flood_reduction	2	4.1	Can this node participate in flood reduction.
hierarchy_indications	3	4.1	Does this node restrict itself to be top-of-fabric or leaf only (in ZTP) and does it support leaf-2-leaf procedures.

## 8.2.19. Registry RIFT\_v4/encoding/NodeFlags

Indication flags of the node.

## 8.2.19.1. Requested Entries

Name	Value	Schema Version	Description
overload	1	4.1	Indicates that node is in overload, do not transit traffic through it.

## 8.2.20. Registry RIFT\_v4/encoding/NodeNeighborsTIEElement

neighbor of a node

## 8.2.20.1. Requested Entries

Name	Value	Schema Version	Description
level	1	4.1	level of neighbor
cost	3	4.1	Cost to neighbor.
link_ids	4	4.1	can carry description of multiple parallel links in a TIE
bandwidth	5	4.1	total bandwidth to neighbor, this will be normally sum of the bandwidths of all the parallel links.

## 8.2.21. Registry RIFT\_v4/encoding/NodeTIEElement

Description of a node.

It may occur multiple times in different TIEs but if either capabilities values do not match or

flags values do not match or

neighbors repeat with different values

the behavior is undefined and a warning SHOULD be generated. Neighbors can be distributed across multiple TIEs however if the sets are disjoint. Miscablings SHOULD be repeated in every node TIE, otherwise the behavior is undefined.

@note: Observe that absence of fields implies defined defaults.

#### 8.2.21.1. Requested Entries

Name	Value	Schema Version	Description
level	1	4.1	Level of the node.
neighbors	2	4.1	Node's neighbors. If neighbor systemID repeats in other node TIEs of same node the behavior is undefined.
capabilities	3	4.1	Capabilities of the node.
flags	4	4.1	Flags of the node.
name	5	4.1	Optional node name for easier operations.
pod	6	4.1	Pod to which the node belongs.
startup_time	7	4.1	optional startup time of the node
miscabled_links	10	4.1	If any local links are miscabled, the indication is flooded.

#### 8.2.22. Registry RIFT\_v4/encoding/PacketContent

Content of a RIFT packet.

##### 8.2.22.1. Requested Entries

Name	Value	Schema Version	Description
lie	1	4.1	
tide	2	4.1	
tire	3	4.1	
tie	4	4.1	

#### 8.2.23. Registry RIFT\_v4/encoding/PacketHeader

Common RIFT packet header.

## 8.2.23.1. Requested Entries

Name	Value	Schema Version	Description
major_version	1	4.1	Major version of protocol.
minor_version	2	4.1	Minor version of protocol.
sender	3	4.1	Node sending the packet, in case of LIE/TIRE/TIDE also the originator of it.
level	4	4.1	Level of the node sending the packet, required on everything except LIEs. Lack of presence on LIEs indicates UNDEFINED_LEVEL and is used in ZTP procedures.

## 8.2.24. Registry RIFT\_v4/encoding/PrefixAttributes

Attributes of a prefix.

## 8.2.24.1. Requested Entries

Name	Value	Schema Version	Description
metric	2	4.1	Distance of the prefix.
tags	3	4.1	Generic unordered set of route tags, can be redistributed to other protocols or use within the context of real time analytics.
monotonic_clock	4	4.1	Monotonic clock for mobile addresses.
loopback	6	4.1	Indicates if the interface is a node loopback.
directly_attached	7	4.1	Indicates that the prefix is directly attached, i.e. should be routed to even if the node is in overload.
from_link	10	4.1	In case of locally originated prefixes, i.e. interface addresses this can describe which link the address belongs to.

## 8.2.25. Registry RIFT\_v4/encoding/PrefixTIEElement

TIE carrying prefixes

## 8.2.25.1. Requested Entries

Name	Value	Schema Version	Description
prefixes	1	4.1	Prefixes with the associated attributes. If the same prefix repeats in multiple TIEs of same node behavior is unspecified.

## 8.2.26. Registry RIFT\_v4/encoding/ProtocolPacket

RIFT packet structure.

## 8.2.26.1. Requested Entries

Name	Value	Schema Version	Description
header	1	4.1	
content	2	4.1	

## 8.2.27. Registry RIFT\_v4/encoding/TIDEPacket

TIDE with sorted TIE headers, if headers are unsorted, behavior is undefined.

## 8.2.27.1. Requested Entries

Name	Value	Schema Version	Description
start_range	1	4.1	First TIE header in the tide packet.
end_range	2	4.1	Last TIE header in the tide packet.
headers	3	4.1	_Sorted_ list of headers.

## 8.2.28. Registry RIFT\_v4/encoding/TIEElement

Single element in a TIE.

Schema enum `common.TIETypeType` in TIEID indicates which elements MUST be present in the TIEElement. In case of mismatch the unexpected elements MUST be ignored. In case of lack of expected element the TIE an error MUST be reported and the TIE MUST be ignored.

This type can be extended with new optional elements for new `common.TIETypeType` values without breaking the major but if it is necessary to understand whether all nodes support the new type a node capability must be added as well.

## 8.2.28.1. Requested Entries

Name	Value	Schema Version	Description
node	1	4.1	Used in case of enum <code>common.TIETypeType.NodeTIEType</code> .
prefixes	2	4.1	Used in case of enum <code>common.TIETypeType.PrefixTIEType</code> .
positive_disaggregation_prefixes	3	4.1	Positive prefixes (always southbound). It MUST NOT be advertised within a North TIE and ignored otherwise.
negative_disaggregation_prefixes	5	4.1	Transitive, negative prefixes (always southbound) which MUST be aggregated and propagated according to the specification southwards towards lower levels to heal pathological upper level partitioning, otherwise blackholes may occur in multiplane fabrics. It MUST NOT be advertised within a North TIE.
external_prefixes	6	4.1	Externally reimported prefixes.
positive_external_disaggregation_prefixes	7	4.1	Positive external disaggregated prefixes (always southbound). It MUST NOT be advertised within a North TIE and ignored otherwise.
keyvalues	9	4.1	Key-Value store elements.

## 8.2.29. Registry RIFT\_v4/encoding/TIEHeader

Header of a TIE.

@note: TIEID space is a total order achieved by comparing the elements in sequence defined and comparing each value as an unsigned integer of according length.

@note: After sequence number the lifetime received on the envelope must be used for comparison before further fields.

@note: 'origination\_time' and 'origination\_lifetime' are normally disregarded for comparison purposes and carried purely for debugging/security purposes if present. They may be used for comparison of last resort to differentiate otherwise equal ties

#### 8.2.29.1. Requested Entries

Name	Value	Schema Version	Description
tieid	2	4.1	ID of the tie.
seq_nr	3	4.1	Sequence number of the tie.
origination_time	10	4.1	Absolute timestamp when the TIE was generated. This can be used on fabrics with synchronized clock to prevent lifetime modification attacks.
origination_lifetime	12	4.1	Original lifetime when the TIE was generated. This can be used on fabrics with synchronized clock to prevent lifetime modification attacks.

#### 8.2.30. Registry RIFT\_v4/encoding/TIEHeaderWithLifeTime

Header of a TIE as described in TIRE/TIDE.

##### 8.2.30.1. Requested Entries

Name	Value	Schema Version	Description
header	1	4.1	
remaining_lifetime	2	4.1	Remaining lifetime that expires down to 0 just like in ISIS. TIEs with lifetimes differing by less than 'lifetime_diff2ignore' MUST be considered EQUAL.

#### 8.2.31. Registry RIFT\_v4/encoding/TIEID

ID of a TIE.

@note: TIEID space is a total order achieved by comparing the elements in sequence defined and comparing each value as an unsigned integer of according length.

## 8.2.31.1. Requested Entries

Name	Value	Schema	Version	Description
direction	1		4.1	direction of TIE
originator	2		4.1	indicates originator of the TIE
tietype	3		4.1	type of the tie
tie_nr	4		4.1	number of the tie

## 8.2.32. Registry RIFT\_v4/encoding/TIEPacket

TIE packet

## 8.2.32.1. Requested Entries

Name	Value	Schema	Version	Description
header	1		4.1	
element	2		4.1	

## 8.2.33. Registry RIFT\_v4/encoding/TIREPacket

TIRE packet

## 8.2.33.1. Requested Entries

Name	Value	Schema	Version	Description
headers	1		4.1	

## 9. Acknowledgments

A new routing protocol in its complexity is not a product of a parent but of a village as the author list shows already. However, many more people provided input, fine-combed the specification based on their experience in design, implementation or application of protocols in IP fabrics. This section will make an inadequate attempt in recording their contribution.

Many thanks to Naiming Shen for some of the early discussions around the topic of using IGPs for routing in topologies related to Clos. Russ White to be especially acknowledged for the key conversation on epistemology that allowed to tie current asynchronous distributed systems theory results to a modern protocol design presented in this scope. Adrian Farrel, Joel Halpern, Jeffrey Zhang, Krzysztof Szarkowicz, Nagendra Kumar, Melchior Aelmans, Kaushal Tank, Will Jones, Moin Ahmed, Sandy Zhang and Jordan Head (in no particular order) provided thoughtful comments that improved the readability of the document and found good amount of corners where the light failed to shine. Kris Price was first to mention single router, single arm default considerations. Jeff Tantsura helped out with some initial

thoughts on BFD interactions while Jeff Haas corrected several misconceptions about BFD's finer points and helped to improve the security section around leaf considerations. Artur Makutunowicz pointed out many possible improvements and acted as sounding board in regard to modern protocol implementation techniques RIFT is exploring. Barak Gafni formalized first time clearly the problem of partitioned spine and fallen leaves on a (clean) napkin in Singapore that led to the very important part of the specification centered around multiple Top-of-Fabric planes and negative disaggregation. Igor Gashinsky and others shared many thoughts on problems encountered in design and operation of large-scale data center fabrics. Xu Benchong found a delicate error in the flooding procedures and a schema datatype size mismatch.

Last but not least, Alvaro Retana guided the undertaking by asking many necessary procedural and technical questions which did not only improve the content but did also lay out the track towards publication.

## 10. References

### 10.1. Normative References

- [EUI64] IEEE, "Guidelines for Use of Extended Unique Identifier (EUI), Organizationally Unique Identifier (OUI), and Company ID (CID)", IEEE EUI, <<http://standards.ieee.org/develop/regauth/tut/eui.pdf>>.
- [ISO10589] ISO "International Organization for Standardization", "Intermediate system to Intermediate system intra-domain routing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473), ISO/IEC 10589:2002, Second Edition.", Nov 2002.
- [RFC1982] Elz, R. and R. Bush, "Serial Number Arithmetic", RFC 1982, DOI 10.17487/RFC1982, August 1996, <<https://www.rfc-editor.org/info/rfc1982>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC2365] Meyer, D., "Administratively Scoped IP Multicast", BCP 23, RFC 2365, DOI 10.17487/RFC2365, July 1998, <<https://www.rfc-editor.org/info/rfc2365>>.



- [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>>.
- [RFC4291] Hinden, R. and S. Deering, "IP Version 6 Addressing Architecture", RFC 4291, DOI 10.17487/RFC4291, February 2006, <<https://www.rfc-editor.org/info/rfc4291>>.
- [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, <<https://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, <<https://www.rfc-editor.org/info/rfc5120>>.
- [RFC5303] Katz, D., Saluja, R., and D. Eastlake 3rd, "Three-Way Handshake for IS-IS Point-to-Point Adjacencies", RFC 5303, DOI 10.17487/RFC5303, October 2008, <<https://www.rfc-editor.org/info/rfc5303>>.
- [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>>.
- [RFC5709] Bhatia, M., Manral, V., Fanto, M., White, R., Barnes, M., Li, T., and R. Atkinson, "OSPFv2 HMAC-SHA Cryptographic Authentication", RFC 5709, DOI 10.17487/RFC5709, October 2009, <<https://www.rfc-editor.org/info/rfc5709>>.
- [RFC5881] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD) for IPv4 and IPv6 (Single Hop)", RFC 5881, DOI 10.17487/RFC5881, June 2010, <<https://www.rfc-editor.org/info/rfc5881>>.
- [RFC5905] Mills, D., Martin, J., Ed., Burbank, J., and W. Kasch, "Network Time Protocol Version 4: Protocol and Algorithms Specification", RFC 5905, DOI 10.17487/RFC5905, June 2010, <<https://www.rfc-editor.org/info/rfc5905>>.

- [RFC6830] Farinacci, D., Fuller, V., Meyer, D., and D. Lewis, "The Locator/ID Separation Protocol (LISP)", RFC 6830, DOI 10.17487/RFC6830, January 2013, <<https://www.rfc-editor.org/info/rfc6830>>.
- [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>>.
- [RFC7987] Ginsberg, L., Wells, P., Decraene, B., Przygienda, T., and H. Gredler, "IS-IS Minimum Remaining Lifetime", RFC 7987, DOI 10.17487/RFC7987, October 2016, <<https://www.rfc-editor.org/info/rfc7987>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8200] Deering, S. and R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification", STD 86, RFC 8200, DOI 10.17487/RFC8200, July 2017, <<https://www.rfc-editor.org/info/rfc8200>>.
- [RFC8202] Ginsberg, L., Previdi, S., and W. Henderickx, "IS-IS Multi-Instance", RFC 8202, DOI 10.17487/RFC8202, June 2017, <<https://www.rfc-editor.org/info/rfc8202>>.
- [RFC8505] Thubert, P., Ed., Nordmark, E., Chakrabarti, S., and C. Perkins, "Registration Extensions for IPv6 over Low-Power Wireless Personal Area Network (6LoWPAN) Neighbor Discovery", RFC 8505, DOI 10.17487/RFC8505, November 2018, <<https://www.rfc-editor.org/info/rfc8505>>.
- [thrift] Apache Software Foundation, "Thrift Interface Description Language", <<https://thrift.apache.org/docs/idl>>.

## 10.2. Informative References

- [CLOS] Yuan, X., "On Nonblocking Folded-Clos Networks in Computer Communication Environments", IEEE International Parallel & Distributed Processing Symposium, 2011.
- [DIJKSTRA] Dijkstra, E., "A Note on Two Problems in Connexion with Graphs", Journal Numer. Math. , 1959.

- [DOT] Ellson, J. and L. Koutsofios, "Graphviz: open source graph drawing tools", Springer-Verlag , 2001.
- [DYNAMO] De Candia et al., G., "Dynamo: amazon's highly available key-value store", ACM SIGOPS symposium on Operating systems principles (SOSP '07), 2007.
- [EPPSTEIN] Eppstein, D., "Finding the k-Shortest Paths", 1997.
- [FATTREE] Leiserson, C., "Fat-Trees: Universal Networks for Hardware-Efficient Supercomputing", 1985.
- [IEEEstd1588] IEEE, "IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems", IEEE Standard 1588, <<https://ieeexplore.ieee.org/document/4579760/>>.
- [IEEEstd8021AS] IEEE, "IEEE Standard for Local and Metropolitan Area Networks - Timing and Synchronization for Time-Sensitive Applications in Bridged Local Area Networks", IEEE Standard 802.1AS, <<https://ieeexplore.ieee.org/document/5741898/>>.
- [ISO10589-Second-Edition] International Organization for Standardization, "Intermediate system to Intermediate system intra-domain routing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473)", Nov 2002.
- [RFC0826] Plummer, D., "An Ethernet Address Resolution Protocol: Or Converting Network Protocol Addresses to 48.bit Ethernet Address for Transmission on Ethernet Hardware", STD 37, RFC 826, DOI 10.17487/RFC0826, November 1982, <<https://www.rfc-editor.org/info/rfc826>>.
- [RFC2131] Droms, R., "Dynamic Host Configuration Protocol", RFC 2131, DOI 10.17487/RFC2131, March 1997, <<https://www.rfc-editor.org/info/rfc2131>>.
- [RFC3626] Clausen, T., Ed. and P. Jacquet, Ed., "Optimized Link State Routing Protocol (OLSR)", RFC 3626, DOI 10.17487/RFC3626, October 2003, <<https://www.rfc-editor.org/info/rfc3626>>.

- [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, <<https://www.rfc-editor.org/info/rfc4861>>.
- [RFC4862] Thomson, S., Narten, T., and T. Jinmei, "IPv6 Stateless Address Autoconfiguration", RFC 4862, DOI 10.17487/RFC4862, September 2007, <<https://www.rfc-editor.org/info/rfc4862>>.
- [RFC6518] Lebovitz, G. and M. Bhatia, "Keying and Authentication for Routing Protocols (KARP) Design Guidelines", RFC 6518, DOI 10.17487/RFC6518, February 2012, <<https://www.rfc-editor.org/info/rfc6518>>.
- [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>>.
- [RFC8415] Mrugalski, T., Siodelski, M., Volz, B., Yourtchenko, A., Richardson, M., Jiang, S., Lemon, T., and T. Winters, "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)", RFC 8415, DOI 10.17487/RFC8415, November 2018, <<https://www.rfc-editor.org/info/rfc8415>>.
- [VAHDAT08] Al-Fares, M., Loukissas, A., and A. Vahdat, "A Scalable, Commodity Data Center Network Architecture", SIGCOMM , 2008.
- [Wikipedia] Wikipedia, "[https://en.wikipedia.org/wiki/Serial\\_number\\_arithmetic](https://en.wikipedia.org/wiki/Serial_number_arithmetic)", 2016.

#### Appendix A. Sequence Number Binary Arithmetic

The only reasonably reference to a cleaner than [RFC1982] sequence number solution is given in [Wikipedia]. It basically converts the problem into two complement's arithmetic. Assuming a straight two complement's subtractions on the bit-width of the sequence number the according >: and =: relations are defined as:

U\_1, U\_2 are 12-bits aligned unsigned version number

D\_f is ( U\_1 - U\_2 ) interpreted as two complement signed 12-bits

D\_b is ( U\_2 - U\_1 ) interpreted as two complement signed 12-bits

U\_1 >: U\_2 IIF D\_f > 0 AND D\_b < 0

U\_1 =: U\_2 IIF D\_f = 0

The >: relationship is anti-symmetric but not transitive. Observe that this leaves >: of the numbers having maximum two complement distance, e.g. ( 0 and 0x800 ) undefined in our 12-bits case since D\_f and D\_b are both -0x7ff.

A simple example of the relationship in case of 3-bit arithmetic follows as table indicating D\_f/D\_b values and then the relationship of U\_1 to U\_2:

U2 / U1	0	1	2	3	4	5	6	7
0	+/+	+/-	+/-	+/-	-/-	-/+	-/+	-/+
1	-/+	+/+	+/-	+/-	+/-	-/-	-/+	-/+
2	-/+	-/+	+/+	+/-	+/-	+/-	-/-	-/+
3	-/+	-/+	-/+	+/+	+/-	+/-	+/-	-/-
4	-/-	-/+	-/+	-/+	+/+	+/-	+/-	+/-
5	+/-	-/-	-/+	-/+	-/+	+/+	+/-	+/-
6	+/-	+/-	-/-	-/+	-/+	-/+	+/+	+/-
7	+/-	+/-	+/-	-/-	-/+	-/+	-/+	+/+

U2 / U1	0	1	2	3	4	5	6	7
0	=	>	>	>	?	<	<	<
1	<	=	>	>	>	?	<	<
2	<	<	=	>	>	>	?	<
3	<	<	<	=	>	>	>	?
4	?	<	<	<	=	>	>	>
5	>	?	<	<	<	=	>	>
6	>	>	?	<	<	<	=	>
7	>	>	>	?	<	<	<	=

#### Appendix B. Information Elements Schema

This section introduces the schema for information elements. The IDL is Thrift [thrift].

On schema changes that

1. change field numbers or
2. add new \*required\* fields or

3. remove any fields or
4. change lists into sets, unions into structures or
5. change multiplicity of fields or
6. changes name of any field or type or
7. change data types of any field or
8. adds, changes or removes a default value of any *\*existing\** field or
9. removes or changes any defined constant or constant value or
10. changes any enumeration type except extending `'common.TIETypeType'` (use of enumeration types is generally discouraged)

major version of the schema MUST increase. All other changes MUST increase minor version within the same major.

The above set of rules guarantees that every decoder can process serialized content generated by a higher minor version of the schema and with that the protocol can progress without a 'fork-lift'. Additionally, based on the propagated minor version in encoded content and added optional node capabilities new TIE types or even de-facto mandatory fields can be introduced without progressing the major version albeit only nodes supporting such new extensions would decode them. Given the model is encoded at the source and never re-encoded flooding through nodes not understanding any new extensions will preserve the according fields.

Content serialized using a major version X is NOT expected to be decodable by any implementation using decoder for a model with a major version lower than X.

Observe especially that introducing an optional field does not cause a major version increase even if the fields inside the structure are optional with defaults.

All signed integer as forced by Thrift [thrift] support must be cast for internal purposes to equivalent unsigned values without discarding the signedness bit. An implementation SHOULD try to avoid using the signedness bit when generating values.

The schema is normative.

## B.1. common.thrift

```
/**
 * Thrift file with common definitions for RIFT
 */

/** @note MUST be interpreted in implementation as unsigned 64 bits.
 *     The implementation SHOULD NOT use the MSB.
 */
typedef i64      SystemIDType
typedef i32      IPv4Address
/** this has to be long enough to accomodate prefix */
typedef binary  IPv6Address
/** @note MUST be interpreted in implementation as unsigned */
typedef i16      UDPPortType
/** @note MUST be interpreted in implementation as unsigned */
typedef i32      TIENrType
/** @note MUST be interpreted in implementation as unsigned */
typedef i32      MTUSizeType
/** @note MUST be interpreted in implementation as unsigned
 *     rolling over number */
typedef i64      SeqNrType
/** @note MUST be interpreted in implementation as unsigned */
typedef i32      LifeTimeInSecType
/** @note MUST be interpreted in implementation as unsigned */
typedef i8       LevelType
/** optional, recommended monotonically increasing number
 *     _per packet type per adjacency_
 *     that can be used to detect losses/misordering/restarts.
 *     @note MUST be interpreted in implementation as unsigned
 *     rolling over number */
typedef i16      PacketNumberType
/** @note MUST be interpreted in implementation as unsigned */
typedef i32      PodType
/** @note MUST be interpreted in implementation as unsigned.
 *     This is carried in the
 *     security envelope and MUST fit into 8 bits. */
typedef i8       VersionType
/** @note MUST be interpreted in implementation as unsigned */
typedef i16      MinorVersionType
/** @note MUST be interpreted in implementation as unsigned */
typedef i32      MetricType
/** @note MUST be interpreted in implementation as unsigned
 *     and unstructured */
typedef i64      RouteTagType
/** @note MUST be interpreted in implementation as unstructured
```

```
        label value */
typedef i32      LabelType
/** @note MUST be interpreted in implementation as unsigned */
typedef i32      BandwidthInMegaBitsType
/** @note Key Value key ID type */
typedef string   KeyIDType
/** node local, unique identification for a link (interface/tunnel
 * etc. Basically anything RIFT runs on). This is kept
 * at 32 bits so it aligns with BFD [RFC5880] discriminator size.
 */
typedef i32      LinkIDType
typedef string   KeyNameType
typedef i8       PrefixLenType
/** timestamp in seconds since the epoch */
typedef i64      TimestampInSecsType
/** security nonce.
 * @note MUST be interpreted in implementation as rolling
 * over unsigned value */
typedef i16      NonceType
/** LIE FSM holdtime type */
typedef i16      TimeIntervalInSecType
/** Transaction ID type for prefix mobility as specified by RFC6550,
 * value MUST be interpreted in implementation as unsigned */
typedef i8       PrefixTransactionIDType
/** Timestamp per IEEE 802.1AS, all values MUST be interpreted in
 * implementation as unsigned. */
struct IEEE802_1ASTimeStampType {
    1: required      i64      AS_sec;
    2: optional      i32      AS_nsec;
}
/** generic counter type */
typedef i64      CounterType
/** Platform Interface Index type, i.e. index of interface on hardware,
 * can be used e.g. with RFC5837 */
typedef i32      PlatformInterfaceIndex

/** Flags indicating node configuration in case of ZTP.
 */
enum HierarchyIndications {
    /** forces level to `leaf_level` and enables according procedures */
    leaf_only = 0,
    /** forces level to `leaf_level` and enables according procedures */
    leaf_only_and_leaf_2_leaf_procedures = 1,
    /** forces level to `top_of_fabric` and enables according
     * procedures */
    top_of_fabric = 2,
}

```



```
const PacketNumberType undefined_packet_number = 0
/** This MUST be used when node is configured as top of fabric in ZTP.
    This is kept reasonably low to allow for fast ZTP convergence on
    failures. */
const LevelType top_of_fabric_level = 24
/** default bandwidth on a link */
const BandwithInMegaBitsType default_bandwidth = 100
/** fixed leaf level when ZTP is not used */
const LevelType leaf_level = 0
const LevelType default_level = leaf_level
const PodType default_pod = 0
const LinkIDType undefined_linkid = 0

/** default distance used */
const MetricType default_distance = 1
/** any distance larger than this will be considered infinity */
const MetricType infinite_distance = 0x7FFFFFFF
/** represents invalid distance */
const MetricType invalid_distance = 0
const bool overload_default = false
const bool flood_reduction_default = true
/** default LIE FSM holddown time */
const TimeIntervalInSecType default_lie_holdtime = 3
/** default ZTP FSM holddown time */
const TimeIntervalInSecType default_ztp_holdtime = 1
/** by default LIE levels are ZTP offers */
const bool default_not_a_ztp_offer = false
/** by default everyone is repeating flooding */
const bool default_you_are_flood_repeater = true
/** 0 is illegal for SystemID */
const SystemIDType IllegalSystemID = 0
/** empty set of nodes */
const set<SystemIDType> empty_set_of_nodeids = {}
/** default lifetime of TIE is one week */
const LifeTimeInSecType default_lifetime = 604800
/** default lifetime when TIEs are purged is 5 minutes */
const LifeTimeInSecType purge_lifetime = 300
/** round down interval when TIEs are sent with security hashes
    to prevent excessive computation. */
const LifeTimeInSecType rounddown_lifetime_interval = 60
/** any 'TieHeader' that has a smaller lifetime difference
    than this constant is equal (if other fields equal). This
    constant MUST be larger than 'purge_lifetime' to avoid
    retransmissions */
const LifeTimeInSecType lifetime_diff2ignore = 400

/** default UDP port to run LIEs on */
const UDPPortType default_lie_udp_port = 914
```

```
/** default UDP port to receive TIEs on, that can be peer specific */
const UDPPortType    default_tie_udp_flood_port = 915

/** default MTU link size to use */
const MTUSizeType    default_mtu_size          = 1400
/** default link being BFD capable */
const bool           bfd_default               = true

/** undefined nonce, equivalent to missing nonce */
const NonceType      undefined_nonce          = 0;
/** outer security key id, MUST be interpreted as in implementation
    as unsigned */
typedef i8           OuterSecurityKeyID
/** security key id, MUST be interpreted as in implementation
    as unsigned */
typedef i32          TIESecurityKeyID
/** undefined key */
const TIESecurityKeyID undefined_securitykey_id = 0;
/** Maximum delta (negative or positive) that a mirrored nonce can
    deviate from local value to be considered valid. If nonces are
    changed every minute on both sides this opens statistically
    a 'maximum_valid_nonce_delta' minutes window of identical LIEs,
    TIE, TI(x)E replays.
    The interval cannot be too small since LIE FSM may change
    states fairly quickly during ZTP without sending LIEs*/
const i16            maximum_valid_nonce_delta = 5;

/** Direction of TIEs. */
enum TieDirectionType {
    Illegal          = 0,
    South             = 1,
    North            = 2,
    DirectionMaxValue = 3,
}

/** Address family type. */
enum AddressFamilyType {
    Illegal          = 0,
    AddressFamilyMinValue = 1,
    IPv4             = 2,
    IPv6             = 3,
    AddressFamilyMaxValue = 4,
}

/** IPv4 prefix type. */
struct IPv4PrefixType {
    1: required IPv4Address    address;
    2: required PrefixLenType  prefixlen;
}
```

```
}

/** IPv6 prefix type. */
struct IPv6PrefixType {
    1: required IPv6Address    address;
    2: required PrefixLenType  prefixlen;
}

/** IP address type. */
union IPAddressType {
    /** Content is IPv4 */
    1: optional IPv4Address    ipv4address;
    /** Content is IPv6 */
    2: optional IPv6Address    ipv6address;
}

/** Prefix advertisement.

    @note: for interface
           addresses the protocol can propagate the address part beyond
           the subnet mask and on reachability computation that has to
           be normalized. The non-significant bits can be used
           for operational purposes.
*/
union IPPrefixType {
    1: optional IPv4PrefixType  ipv4prefix;
    2: optional IPv6PrefixType  ipv6prefix;
}

/** Sequence of a prefix in case of move.
*/
struct PrefixSequenceType {
    1: required IEEE802_1ASTimeStampType  timestamp;
    /** Transaction ID set by client in e.g. in 6LoWPAN. */
    2: optional PrefixTransactionIDType  transactionid;
}

/** Type of TIE.

    This enum indicates what TIE type the TIE is carrying.
    In case the value is not known to the receiver,
    the TIE MUST be re-flooded. This allows for
    future extensions of the protocol within the same major schema
    with types opaque to some nodes UNLESS the flooding scope is not
    the same as prefix TIE, then a major version revision MUST
    be performed.
*/
enum TIETypeType {
```

```

    Illegal                = 0,
    TIETypeMinValue        = 1,
    /** first legal value */
    NodeTIEType            = 2,
    PrefixTIEType          = 3,
    PositiveDisaggregationPrefixTIEType = 4,
    NegativeDisaggregationPrefixTIEType = 5,
    PGPrefixTIEType        = 6,
    KeyValueTIEType        = 7,
    ExternalPrefixTIEType  = 8,
    PositiveExternalDisaggregationPrefixTIEType = 9,
    TIETypeMaxValue        = 10,
}

```

```
/** RIFT route types.
```

```

    @note: route types which MUST be ordered on their preference
           PGP prefixes are most preferred attracting
           traffic north (towards spine) and then south
           normal prefixes are attracting traffic south
           (towards leafs), i.e. prefix in NORTH PREFIX TIE
           is preferred over SOUTH PREFIX TIE.

```

```

    @note: The only purpose of those values is to introduce an
           ordering whereas an implementation can choose internally
           any other values as long the ordering is preserved

```

```

*/
enum RouteType {
    Illegal                = 0,
    RouteTypeMinValue      = 1,
    /** First legal value. */
    /** Discard routes are most preferred */
    Discard                 = 2,

    /** Local prefixes are directly attached prefixes on the
     * system such as e.g. interface routes.
     */
    LocalPrefix             = 3,
    /** Advertised in S-TIEs */
    SouthPGPPrefix          = 4,
    /** Advertised in N-TIEs */
    NorthPGPPrefix          = 5,
    /** Advertised in N-TIEs */
    NorthPrefix             = 6,
    /** Externally imported north */
    NorthExternalPrefix     = 7,
    /** Advertised in S-TIEs, either normal prefix or positive
     * disaggregation */

```

```
    SouthPrefix          = 8,  
    /** Externally imported south */  
    SouthExternalPrefix = 9,  
    /** Negative, transitive prefixes are least preferred */  
    NegativeSouthPrefix = 10,  
    RouteTypeMaxValue   = 11,  
}
```

## B.2. encoding.thrift

```
/**  
    Thrift file for packet encodings for RIFT  
*/  
  
include "common.thrift"  
  
namespace rs models  
namespace py encoding  
  
/** Represents protocol encoding schema major version */  
const common.VersionType protocol_major_version = 4  
/** Represents protocol encoding schema minor version */  
const common.MinorVersionType protocol_minor_version = 1  
  
/** Common RIFT packet header. */  
struct PacketHeader {  
    /** Major version of protocol. */  
    1: required common.VersionType    major_version =  
        protocol_major_version;  
    /** Minor version of protocol. */  
    2: required common.MinorVersionType minor_version =  
        protocol_minor_version;  
    /** Node sending the packet, in case of LIE/TIRE/TIDE  
        also the originator of it. */  
    3: required common.SystemIDType  sender;  
    /** Level of the node sending the packet, required on everything  
        except LIEs. Lack of presence on LIEs indicates UNDEFINED_LEVEL  
        and is used in ZTP procedures.  
    */  
    4: optional common.LevelType      level;  
}  
  
/** Prefix community. */  
struct Community {
```

```
    /** Higher order bits */
    1: required i32          top;
    /** Lower order bits */
    2: required i32          bottom;
}

/** Neighbor structure. */
struct Neighbor {
    /** System ID of the originator. */
    1: required common.SystemIDType    originator;
    /** ID of remote side of the link. */
    2: required common.LinkIDType      remote_id;
}

/** Capabilities the node supports.

@note: The schema may add to this
field future capabilities to indicate whether it will support
interpretation of future schema extensions on the same major
revision. Such fields MUST be optional and have an implicit or
explicit false default value. If a future capability changes route
selection or generates blackholes if some nodes are not supporting
it then a major version increment is unavoidable.
*/
struct NodeCapabilities {
    /** Must advertise supported minor version dialect that way. */
    1: required common.MinorVersionType    protocol_minor_version =
        protocol_minor_version;
    /** Can this node participate in flood reduction. */
    2: optional bool                       flood_reduction =
        common.flood_reduction_default;
    /** Does this node restrict itself to be top-of-fabric or
leaf only (in ZTP) and does it support leaf-2-leaf
procedures. */
    3: optional common.HierarchyIndications    hierarchy_indications;
}

/** Link capabilities. */
struct LinkCapabilities {
    /** Indicates that the link is supporting BFD. */
    1: optional bool                       bfd =
        common.bfd_default;
    /** Indicates whether the interface will support v4 forwarding.

@note: This MUST be set to true when LIEs from a v4 address are
sent and MAY be set to true in LIEs on v6 address. If v4
and v6 LIEs indicate contradicting information the
behavior is unspecified. */

```

```
        2: optional bool                                v4_forwarding_capable =
            true;
    }

/** RIFT LIE Packet.

    @note: this node's level is already included on the packet header
*/
struct LIEPacket {
    /** Node or adjacency name. */
    1: optional string                                name;
    /** Local link ID. */
    2: required common.LinkIDType                    local_id;
    /** UDP port to which we can receive flooded TIEs. */
    3: required common.UDPPortType                  flood_port =
        common.default_tie_udp_flood_port;
    /** Layer 3 MTU, used to discover to mismatch. */
    4: optional common.MTUSizeType                  link_mtu_size =
        common.default_mtu_size;
    /** Local link bandwidth on the interface. */
    5: optional common.BandwidthInMegaBitsType      link_bandwidth = common.default_bandwidth;
    /** Reflects the neighbor once received to provide
        3-way connectivity. */
    6: optional Neighbor                             neighbor;
    /** Node's PoD. */
    7: optional common.PodType                       pod =
        common.default_pod;
    /** Node capabilities shown in LIE. The capabilities
        MUST match the capabilities shown in the Node TIEs, otherwise
        the behavior is unspecified. A node detecting the mismatch
        SHOULD generate according error. */
    10: required NodeCapabilities                    node_capabilities;
    /** Capabilities of this link. */
    11: optional LinkCapabilities                   link_capabilities;
    /** Required holdtime of the adjacency, i.e. how much time
        MUST expire without LIE for the adjacency to drop. */
    12: required common.TimeIntervalInSecType       holdtime = common.default_lie_holdtime;
    /** Unsolicited, downstream assigned locally significant label
        value for the adjacency. */
    13: optional common.LabelType                   label;
    /** Indicates that the level on the LIE MUST NOT be used
        to derive a ZTP level by the receiving node. */
    21: optional bool                                not_a_ztp_offer =
        common.default_not_a_ztp_offer;
    /** Indicates to northbound neighbor that it should
        be refloding this node's N-TIEs to achieve flood reduction and
```

```
    balancing for northbound flooding. To be ignored if received
    from a northbound adjacency. */
22: optional bool                you_are_flood_repeater =
    common.default_you_are_flood_repeater;
/** Can be optionally set to indicate to neighbor that packet losses
    are seen on reception based on packet numbers or the rate is
    too high. The receiver SHOULD temporarily slow down
    flooding rates.
    */
23: optional bool                you_are_sending_too_quickly =
    false;
/** Instance name in case multiple RIFT instances running on same
    interface. */
24: optional string              instance_name;
}

/** LinkID pair describes one of parallel links between two nodes. */
struct LinkIDPair {
    /** Node-wide unique value for the local link. */
    1: required common.LinkIDType    local_id;
    /** Received remote link ID for this link. */
    2: required common.LinkIDType    remote_id;

    /** Describes the local interface index of the link. */
    10: optional common.PlatformInterfaceIndex platform_interface_index;
    /** Describes the local interface name. */
    11: optional string              platform_interface_name;
    /** Indication whether the link is secured, i.e. protected by
        outer key, absence of this element means no indication,
        undefined outer key means not secured. */
    12: optional common.OuterSecurityKeyID
        trusted_outer_security_key;
    /** Indication whether the link is protected by established
        BFD session. */
    13: optional bool                bfd_up;
    /** Optional indication which address families are up on the
        interface */
    14: optional set<common.AddressFamilyType> address_families;
}

/** ID of a TIE.

    @note: TIEID space is a total order achieved by comparing
    the elements in sequence defined and comparing each
    value as an unsigned integer of according length.
    */
struct TIEID {
    /** direction of TIE */
```



```
    1: required common.TieDirectionType    direction;
    /** indicates originator of the TIE */
    2: required common.SystemIDType        originator;
    /** type of the tie */
    3: required common.TIETypeType        tietype;
    /** number of the tie */
    4: required common.TIENrType          tie_nr;
}

/** Header of a TIE.

@note: TIEID space is a total order achieved by comparing
the elements in sequence defined and comparing each
value as an unsigned integer of according length.

@note: After sequence number the lifetime received on the envelope
must be used for comparison before further fields.

@note: `origination_time` and `origination_lifetime` are
normally disregarded for comparison purposes and carried
purely for debugging/security purposes if present.
They may be used for comparison of last resort to
differentiate otherwise equal ties
*/
struct TIEHeader {
    /** ID of the tie. */
    2: required TIEID                                tieid;
    /** Sequence number of the tie. */
    3: required common.SeqNrType                    seq_nr;

    /** Absolute timestamp when the TIE
was generated. This can be used on fabrics with
synchronized clock to prevent lifetime modification attacks. */
    10: optional common.IEEE802_1ASTimeStampType    origination_time;
    /** Original lifetime when the TIE
was generated. This can be used on fabrics with
synchronized clock to prevent lifetime modification attacks. */
    12: optional common.LifeTimeInSecType          origination_lifetime;
}

/** Header of a TIE as described in TIRE/TIDE.
*/
struct TIEHeaderWithLifeTime {
    1: required TIEHeader                            header;
    /** Remaining lifetime that expires down to 0 just like in ISIS.
TIEs with lifetimes differing by less than
`lifetime_diff2ignore` MUST be considered EQUAL. */
    2: required common.LifeTimeInSecType            remaining_lifetime;
}
```

```
}

/** TIDE with sorted TIE headers, if headers are unsorted, behavior
    is undefined. */
struct TIDEPacket {
    /** First TIE header in the tide packet. */
    1: required TIEID                start_range;
    /** Last TIE header in the tide packet. */
    2: required TIEID                end_range;
    /** _Sorted_ list of headers. */
    3: required list<TIEHeaderWithLifeTime> headers;
}

/** TIRE packet */
struct TIREPacket {
    1: required set<TIEHeaderWithLifeTime> headers;
}

/** neighbor of a node */
struct NodeNeighborsTIEElement {
    /** level of neighbor */
    1: required common.LevelType      level;
    /** Cost to neighbor.

        @note: All parallel links to same node
        incur same cost, in case the neighbor has multiple
        parallel links at different cost, the largest distance
        (highest numerical value) MUST be advertised.

        @note: any neighbor with cost <= 0 MUST be ignored
        in computations */
    3: optional common.MetricType      cost
        = common.default_distance;
    /** can carry description of multiple parallel links in a TIE */
    4: optional set<LinkIDPair>       link_ids;

    /** total bandwidth to neighbor, this will be normally sum of the
        bandwidths of all the parallel links. */
    5: optional common.BandwidthInMegaBitsType
        bandwidth = common.default_bandwidth;
}

/** Indication flags of the node. */
struct NodeFlags {
    /** Indicates that node is in overload, do not transit traffic
        through it. */
    1: optional bool                  overload = common.overload_default;
}

```

/\*\* Description of a node.

It may occur multiple times in different TIEs but if either

- <t>capabilities values do not match or</t>
- <t>flags values do not match or</t>
- <t>neighbors repeat with different values</t>

the behavior is undefined and a warning SHOULD be generated. Neighbors can be distributed across multiple TIEs however if the sets are disjoint. Miscablings SHOULD be repeated in every node TIE, otherwise the behavior is undefined.

@note: Observe that absence of fields implies defined defaults.

```

*/
struct NodeTIEElement {
    /** Level of the node. */
    1: required common.LevelType          level;
    /** Node's neighbors. If neighbor systemID repeats in other
        node TIEs of same node the behavior is undefined. */
    2: required map<common.SystemIDType,
        NodeNeighborsTIEElement>         neighbors;
    /** Capabilities of the node. */
    3: required NodeCapabilities          capabilities;
    /** Flags of the node. */
    4: optional NodeFlags                  flags;
    /** Optional node name for easier operations. */
    5: optional string                     name;
    /** PoD to which the node belongs. */
    6: optional common.PodType             pod;
    /** optional startup time of the node */
    7: optional common.TimestampInSecsType startup_time;

    /** If any local links are miscabled, the indication is flooded. */
    10: optional set<common.LinkIDType>    miscabled_links;
}

/** Attributes of a prefix. */
struct PrefixAttributes {
    /** Distance of the prefix. */
    2: required common.MetricType          metric
        = common.default_distance;
    /** Generic unordered set of route tags, can be redistributed
        to other protocols or use within the context of real time
        analytics. */
    3: optional set<common.RouteTagType>    tags;
}

```

```
    /** Monotonic clock for mobile addresses. */
    4: optional common.PrefixSequenceType    monotonic_clock;
    /** Indicates if the interface is a node loopback. */
    6: optional bool                        loopback = false;
    /** Indicates that the prefix is directly attached, i.e. should be
        routed to even if the node is in overload. */
    7: optional bool                        directly_attached = true;

    /** In case of locally originated prefixes, i.e. interface
        addresses this can describe which link the address
        belongs to. */
    10: optional common.LinkIDType          from_link;
}

/** TIE carrying prefixes */
struct PrefixTIEElement {
    /** Prefixes with the associated attributes.
        If the same prefix repeats in multiple TIEs of same node
        behavior is unspecified. */
    1: required map<common.IPPrefixType, PrefixAttributes> prefixes;
}

/** Generic key value pairs. */
struct KeyValueTIEElement {
    /** @note: if the same key repeats in multiple TIEs of same node
        or with different values, behavior is unspecified */
    1: required map<common.KeyIDType, string>    keyvalues;
}

/** Single element in a TIE.

    Schema enum `common.TIETypeType`
    in TIEID indicates which elements MUST be present
    in the TIEElement. In case of mismatch the unexpected
    elements MUST be ignored. In case of lack of expected
    element the TIE an error MUST be reported and the TIE
    MUST be ignored.

    This type can be extended with new optional elements
    for new `common.TIETypeType` values without breaking
    the major but if it is necessary to understand whether
    all nodes support the new type a node capability must
    be added as well.
    */
union TIEElement {
    /** Used in case of enum common.TIETypeType.NodeTIEType. */
    1: optional NodeTIEElement    node;
    /** Used in case of enum common.TIETypeType.PrefixTIEType. */

```

```
2: optional PrefixTIEElement    prefixes;
/** Positive prefixes (always southbound).
    It MUST NOT be advertised within a North TIE and
    ignored otherwise.
*/
3: optional PrefixTIEElement    positive_disaggregation_prefixes;
/** Transitive, negative prefixes (always southbound) which
    MUST be aggregated and propagated
    according to the specification
    southwards towards lower levels to heal
    pathological upper level partitioning, otherwise
    blackholes may occur in multiplane fabrics.
    It MUST NOT be advertised within a North TIE.
*/
5: optional PrefixTIEElement    negative_disaggregation_prefixes;
/** Externally reimported prefixes. */
6: optional PrefixTIEElement    external_prefixes;
/** Positive external disaggregated prefixes (always southbound).
    It MUST NOT be advertised within a North TIE and
    ignored otherwise.
*/
7: optional PrefixTIEElement    positive_external_disaggregation_prefixes;
/** Key-Value store elements. */
9: optional KeyValueTIEElement  keyvalues;
}

/** TIE packet */
struct TIEPacket {
    1: required TIEHeader  header;
    2: required TIEElement element;
}

/** Content of a RIFT packet. */
union PacketContent {
    1: optional LIEPacket    lie;
    2: optional TIDEPacket   tide;
    3: optional TIREPacket   tire;
    4: optional TIEPacket    tie;
}

/** RIFT packet structure. */
struct ProtocolPacket {
    1: required PacketHeader header;
    2: required PacketContent content;
}
```

## Appendix C. Constants

### C.1. Configurable Protocol Constants

This section gathers constants that are provided in the schema files and in the document.

	Type	Value
LIE IPv4 Multicast Address	Default Value, Configurable	224.0.0.120 or all-rift-routers to be assigned in IPv4 Multicast Address Space Registry in Local Network Control Block
LIE IPv6 Multicast Address	Default Value, Configurable	FF02::A1F7 or all-rift-routers to be assigned in IPv6 Multicast Address Assignments
LIE Destination Port	Default Value, Configurable	914
Level value for TOP_OF_FABRIC flag	Constant	24
Default LIE Holdtime	Default Value, Configurable	3 seconds
TIE Retransmission Interval	Default Value	1 second
TIDE Generation Interval	Default Value, Configurable	5 seconds
MIN_TIEID signifies start of TIDEs	Constant	TIE Key with minimal values: TIEID(originator=0, tietype=TIETypeMinValue, tie_nr=0, direction=South)
MAX_TIEID signifies end of TIDEs	Constant	TIE Key with maximal values: TIEID(originator=MAX_UINT64, tietype=TIETypeMaxValue, tie_nr=MAX_UINT64, direction=North)

Table 6: all\_constants

Authors' Addresses

Tony Przygienda (editor)  
Juniper  
1137 Innovation Way

Sunnyvale, CA

USA

Email: [prz@juniper.net](mailto:prz@juniper.net)

Alankar Sharma  
Comcast  
1800 Bishops Gate Blvd  
Mount Laurel, NJ 08054  
US

Email: [Alankar\\_Sharma@comcast.com](mailto:Alankar_Sharma@comcast.com)

Pascal Thubert  
Cisco Systems, Inc  
Building D  
45 Allee des Ormes - BP1200  
MOUGINS - Sophia Antipolis 06254  
FRANCE

Phone: +33 497 23 26 34  
Email: [pthubert@cisco.com](mailto:pthubert@cisco.com)

Bruno Rijsman  
Individual

Email: [brunorijsman@gmail.com](mailto:brunorijsman@gmail.com)

Dmitry Afanasiev  
Yandex

Email: [f10w@yandex-team.ru](mailto:f10w@yandex-team.ru)



RIFT WG  
Internet-Draft  
Intended status: Standards Track  
Expires: November 5, 2019

Zheng. Zhang  
Yuehua. Wei  
ZTE Corporation  
Shaowen. Ma  
Mellanox  
Xufeng. Liu  
Volta Networks  
May 4, 2019

RIFT YANG Model  
draft-zhang-rift-yang-02

Abstract

This document defines a YANG data model for the configuration and management of RIFT Protocol.

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 November 5, 2019.

Copyright Notice

Copyright (c) 2019 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 Simplified BSD License text as described in Section 4.e of

the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

1. Introduction . . . . .	2
2. Design of the Data Model . . . . .	2
3. RIFT configuration . . . . .	7
4. RIFT State . . . . .	7
5. RPC . . . . .	7
6. Notifications . . . . .	7
7. RIFT YANG model . . . . .	7
8. Security Considerations . . . . .	19
9. IANA Considerations . . . . .	20
10. Contributors . . . . .	20
11. Normative References . . . . .	20
Authors' Addresses . . . . .	22

## 1. Introduction

[I-D.ietf-rift-rift] introduces the protocol definition of RIFT. This document defines a YANG data model that can be used to configure and manage the RIFT protocol. The model is based on YANG 1.1 as defined in [RFC7950] and conforms to the Network Management Datastore Architecture (NDMA) as described in [RFC8342]

## 2. Design of the Data Model

This model imports and augments ietf-routing YANG model defined in [RFC8349]. Both configuration branch and state branch of [RFC8349] are augmented. The configuration branch covers node base and policy configuration. The neighbor state will be added in later version. The container "rift" is the top level container in this data model. The presence of this container is expected to enable RIFT protocol functionality.

module: ietf-rift

```
augment /rt:routing/rt:control-plane-protocols/rt:control-plane-protocol:
  +--rw rift!
    +--rw node-info
      |
      | +--rw level?                level
      | +--rw systemid            systemid
      | +--rw name?               string
      | +--rw pod?                uint32
      | +--rw overload?           boolean
      | +--rw flood-reducing?     boolean {flood-reducing}?
      | +--rw hierarchy-indications? enumeration
      | +--rw interfaces* [local-id]
```

```

    +--rw local-id                linkidtype
    +--rw name?                   if:interface-ref
    +--rw if-index?               if:interface-ref
    +--rw bfd?                    boolean {bfd}?
    +--rw direction-type?        enumeration
    +--rw you_are_flood_repeater? boolean
    +--rw not_a_ztp_offer?       boolean
    +--rw flood-port?            inet:port-number
    +--rw lie-rx-port?           inet:port-number
    +--rw holdtime?              rt-types:timer-value-seconds16
    +--rw local-nonce?           uint16
+--rw (algorithn-type)?
  |
  |   +--:(spf)
  |   +--:(all-path)
+--ro hal?                       level
+--ro vol-list
  +--ro vol* [systemid]
    +--ro offered-level?         level
    +--ro level?                 level
    +--ro systemid               systemid
    +--ro name?                  string
    +--ro pod?                   uint32
    +--ro overload?              boolean
    +--ro flood-reducing?        boolean {flood-reducing}?
    +--ro hierarchy-indications? enumeration
    +--ro interfaces* [local-id]
      +--ro local-id             linkidtype
      +--ro name?                if:interface-ref
      +--ro if-index?            if:interface-ref
      +--ro bfd?                 boolean {bfd}?
      +--ro direction-type?      enumeration
      +--ro you_are_flood_repeater? boolean
      +--ro not_a_ztp_offer?     boolean
      +--ro flood-port?          inet:port-number
      +--ro lie-rx-port?         inet:port-number
      +--ro holdtime?            rt-types:timer-value-seconds1
  +--ro local-nonce?             uint16
+--ro miscabled-links*          linkidtype
+--ro neighbor
  +--ro nbrs* [systemid remote-id]
    +--ro level?                 level
    +--ro systemid               systemid
    +--ro name?                  string
    +--ro pod?                   uint32
    +--ro overload?              boolean
    +--ro flood-reducing?        boolean {flood-reducing}?
    +--ro hierarchy-indications? enumeration
    +--ro interfaces* [local-id]

```

6



```

    |
    |
    | |
    | | +--ro metric?           uint32
    | | +--ro tag?             uint64
    | | +--ro monotonic_clock? PrefixSequenceType
    | | +--ro from-link?      linkidtype
    | | +--ro positive_disaggregation_prefixes
    | | |
    | | | +--ro prefix?       inet:ip-prefix
    | | | +--ro metric?       uint32
    | | | +--ro tag?          uint64
    | | | +--ro monotonic_clock? PrefixSequenceType
    | | | +--ro from-link?    linkidtype
    | | | +--ro negative_disaggregation_prefixes
    | | | |
    | | | | +--ro prefix?     inet:ip-prefix
    | | | | +--ro metric?     uint32
    | | | | +--ro tag?        uint64
    | | | | +--ro monotonic_clock? PrefixSequenceType
    | | | | +--ro from-link?  linkidtype
    | | | | +--ro external_prefixes
    | | | | |
    | | | | | +--ro prefix?   inet:ip-prefix
    | | | | | +--ro metric?   uint32
    | | | | | +--ro tag?     uint64
    | | | | | +--ro monotonic_clock? PrefixSequenceType
    | | | | | +--ro from-link? linkidtype
    | | | +--ro kvs
    | | | |
    | | | | +--ro key?      uint16
    | | | | +--ro value?   uint32
    | | +--rw kv-store
    | | |
    | | | +--rw kvs* [kvs-index]
    | | | |
    | | | | +--rw kvs-index  uint32
    | | | | +--rw kvs-tie
    | | | | |
    | | | | | +--rw originator?  systemid
    | | | | | +--rw direction-type? enumeration
    | | | | | +--rw tie-number?  uint32
    | | | | | +--rw key?        uint16
    | | | | | +--rw value?     uint32

notifications:
+---n error-set
+--ro tie-level-error
|
| |
| | +--ro originator?      systemid
| | +--ro direction-type? enumeration
| | +--ro tie-number?     uint32
| | +--ro seq?            uint32
| | +--ro lifetime?       uint16
| | +--ro tie-node
| | |
| | | +--ro level?         level
| | | +--ro systemid      systemid
| | | +--ro name?         string
| | | +--ro pod?          uint32

```

```

+--ro overload?                boolean
+--ro flood-reducing?          boolean {flood-reducing}?
+--ro hierarchy-indications?   enumeration
+--ro interfaces* [local-id]
  +--ro local-id                linkidtype
  +--ro name?                   if:interface-ref
  +--ro if-index?              if:interface-ref
  +--ro bfd?                    boolean {bfd}?
  +--ro direction-type?        enumeration
  +--ro you_are_flood_repeater? boolean
  +--ro not_a_ztp_offer?        boolean
  +--ro flood-port?            inet:port-number
  +--ro lie-rx-port?           inet:port-number
  +--ro holdtime?              rt-types:timer-value-seconds16
  +--ro local-nonce?           uint16
+--ro tie-prefix
  +--ro prefixes
    +--ro prefix?              inet:ip-prefix
    +--ro metric?              uint32
    +--ro tag?                 uint64
    +--ro monotonic_clock?     PrefixSequenceType
    +--ro from-link?           linkidtype
  +--ro positive_disaggregation_prefixes
    +--ro prefix?              inet:ip-prefix
    +--ro metric?              uint32
    +--ro tag?                 uint64
    +--ro monotonic_clock?     PrefixSequenceType
    +--ro from-link?           linkidtype
  +--ro negative_disaggregation_prefixes
    +--ro prefix?              inet:ip-prefix
    +--ro metric?              uint32
    +--ro tag?                 uint64
    +--ro monotonic_clock?     PrefixSequenceType
    +--ro from-link?           linkidtype
  +--ro external_prefixes
    +--ro prefix?              inet:ip-prefix
    +--ro metric?              uint32
    +--ro tag?                 uint64
    +--ro monotonic_clock?     PrefixSequenceType
    +--ro from-link?           linkidtype
+--ro kvs
  +--ro key?                    uint16
  +--ro value?                  uint32
+--ro nbr-error
  +--ro nbrs* [systemid remote-id]
    +--ro level?                level
    +--ro systemid              systemid
    +--ro name?                 string

```

```

+--ro pod?                               uint32
+--ro overload?                           boolean
+--ro flood-reducing?                      boolean {flood-reducing}?
+--ro hierarchy-indications?              enumeration
+--ro interfaces* [local-id]
|   +--ro local-id                         linkidtype
|   +--ro name?                            if:interface-ref
|   +--ro if-index?                        if:interface-ref
|   +--ro bfd?                             boolean {bfd}?
|   +--ro direction-type?                  enumeration
|   +--ro you_are_flood_repeater?          boolean
|   +--ro not_a_ztp_offer?                 boolean
|   +--ro flood-port?                      inet:port-number
|   +--ro lie-rx-port?                     inet:port-number
|   +--ro holdtime?                        rt-types:timer-value-seconds16
|   +--ro local-nonce?                     uint16
+--ro remote-id                           uint32
+--ro local-id?                            uint32
+--ro distance?                            uint32
+--ro remote-nonce?                        uint16
+--ro miscabled-links*                     linkidtype

```

### 3. RIFT configuration

RIFT configurations require node base information configurations. Some features can be used to enhance protocol, such as BFD, flooding-reducing, community attribute.

### 4. RIFT State

RIFT states are composed of RIFT node state, neighbor state, database.

### 5. RPC

TBD.

### 6. Notifications

Unexpected TIE and neighbor's layer error should be notified.

### 7. RIFT YANG model

```

<CODE BEGINS> file "ietf-rift.yang"
module ietf-rift {

    yang-version 1.1;

```

```
namespace "urn:ietf:params:xml:ns:yang:ietf-rift";
prefix rift;

import ietf-inet-types {
  prefix "inet";
  reference "RFC6991";
}

import ietf-routing {
  prefix "rt";
  reference "RFC8349";
}

import ietf-interfaces {
  prefix "if";
  reference "RFC7223";
}

import ietf-routing-types {
  prefix "rt-types";
  reference "RFC8294";
}

organization
  "IETF RIFT(Routing In Fat Trees) Working Group";

contact
  "WG Web: <http://tools.ietf.org/wg/rift/>
  WG List: <mailto:rift@ietf.org>

  Editor: Zheng Zhang
  <mailto:zzhang_ietf@hotmail.com>

  Editor: Yuehua Wei
  <mailto:wei.yuehua@zte.com.cn>

  Editor: Shaowen Ma
  <mailto:mashaowen@gmail.com>

  Editor: Xufeng Liu
  <mailto:Xufeng_Liu@jabil.com>";

description
  "The module defines the YANG definitions for RIFT.

  Copyright (c) 2018 IETF Trust and the persons
  identified as authors of the code. All rights reserved.
```



Redistribution and use in source and binary forms, with or without modification, is permitted pursuant to, and subject to the license terms contained in, the Simplified BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>). This version of this YANG module is part of RFC 3618; see the RFC itself for full legal notices.";

```
revision 2019-05-05 {
  description "Initial revision.";
  reference
    "RFC XXXX: A YANG Data Model for RIFT.
     draft-ietf-rift-rift: RIFT: Routing in Fat Trees.";
}

/*
 * Features
 */

/*feature overload {
  description "A node with overload bit set SHOULD NOT advertise any reach
ability
  prefixes southbound except locally hosted ones. The leaf no
de SHOULD
  set the 'overload' bit on its node TIEs.";
}*/

feature bfd {
  description "Support BFD (RFC5881) function to react quickly to link fai
lures.";
}

feature flood-reducing {
  description "Support flood reducing function defined in section 4.2.3.8.
";
}

feature policy {
  description "Support policy guide information.";
}

typedef systemid {
  type string {
    pattern
      '[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}';
  }
  description
    "This type defines RIFT system id using pattern,
    system id looks like : 0143.0438.0100.AeF0";
}
```

```
typedef level {
    type uint8 {
        range "0 .. 24";
    }
    default "0";
    description "The value of node level. The max value is 24.";
}

typedef linkidtype {
    type uint32;
    description
        "This type defines the link id of an interface.";
}

typedef PrefixSequenceType {
    type uint64;
    description
        "This type defines the link id of an interface.";
}

/*
 * Identity
 */
identity rift {
    base rt:routing-protocol;
    description "Identity for the RIFT routing protocol.";
}

/*
 * Groupings
 */
grouping node-key {
    leaf systemid {
        type systemid;
        mandatory true;
        description "Each node is identified via a SystemID which is 64 bits
wide.";
    }
    description "The key information used to distinguish a node.";
}

grouping base-node-info {
    leaf level {
        type level;
        description "The level of this node.";
    }
    uses node-key;

    leaf name {
```

```
        type string;
        description "The name of this node. It won't be used as the key of n
ode,
                                just used for description.";
    }
    leaf pod {
        type uint32;
        description "Point of Delivery. The self-contained vertical slice of
a Clos or Fat Tree network containing normally only
level 0 and level 1 nodes. It communicates with nodes
in other PoDs via the spine. We number PoDs to distinguish
ish
                                them and use PoD #0 to denote 'undefined' PoD.";
    }
    leaf overload {
        type boolean;
        description "If the overload bit in TIEs should be set.";
    }
    leaf flood-reducing {
        if-feature flood-reducing;
        type boolean;
        description "If the node support flood reducing function defined in
section 4.2.3.8.";
    }
    uses hierarchy-indications;
    uses interface;

    description "The base information of a node.";
}

grouping neighbor {
    leaf remote-id {
        type uint32;
        description "The remote-id to reach this neighbor.";
    }
    leaf local-id {
        type uint32;
        description "The local-id of link connect to this neighbor.";
    }
    leaf distance {
        type uint32;
        description "The cost value to arrive this neighbor.";
    }
    leaf remote-nonce {
        type uint16;
        description "Remote Nonce of the adjacency as received
in LIEs. In case of LIE packet this MUST
correspond to the value in the serialized
object otherwise the packet MUST be discarded.";
    }
    leaf-list miscabled-links {
```

```
        type linkidtype;
        description "List of miscabled links.";
    }
    description "The neighbor information.";
}

grouping hierarchy-indications {
    leaf hierarchy-indications {
        type enumeration {
            enum "leaf-only" {
                description "The node will never leave the 'bottom of the
                    hierarchy'.";
            }
            enum "leaf_only_and_leaf_2_leaf_procedures" {
                description "This means leaf to leaf.";
            }
            enum "top_of_fabric" {
                description "The node is 'top of fabric'.";
            }
        }
        description "The hierarchy indications of this node.";
    }
    description "The hierarchy indications of this node.";
}

grouping node {
    uses base-node-info;
    uses algorithm;

    leaf hal {
        type level;
        config false;
        description "The highest defined level value seen from all
            valid level offers received.";
    }
    container vol-list {
        config false;
        list vol {
            key "systemid";
            leaf offered-level {
                type level;
                description "The level type value offered by this neighbor."
            }
        }
        uses base-node-info;
        description "The valid offered level information.";
    }
    description "The valid offered level information.";
}
;
```

```
    leaf-list miscabled-links {
      type linkidtype;
      config false;
      description
        "List of miscabled links.";
    }
    description "The information of local node. Includes base information,
      configurable parameters and features.";
  }

grouping direction-type {
  leaf direction-type {
    type enumeration {
      enum illegal {
        description "Illegal direction.";
      }
      enum south {
        description "A link to a node one level down.";
      }
      enum north {
        description "A link to a node one level up.";
      }
      enum east-west {
        description "A link to a node in the same level.";
      }
      enum max {
        description "The max value of direction.";
      }
    }
    description "The type of a link.";
  }
  description "The type of a link.";
}

grouping interface {
  list interfaces {
    key "local-id";
    leaf local-id {
      type linkidtype;
      mandatory true;
      description "The local id of this interface.";
    }
    leaf name {
      type if:interface-ref;
      description "The interface's name.";
    }
  }
  leaf if-index {
    type if:interface-ref;
  }
}
```

```
        description "The index of this interface.";
    }
    leaf bfd {
        if-feature bfd;
        type boolean;
        description "If BFD function is enabled to react link failures
                    after neighbor's detection.";
    }
    uses direction-type;

    leaf you_are_flood_repeater {
        type boolean;
        description "If the neighbor on this link is flooding repeater."
;
    }
    leaf not_a_ztp_offer {
        type boolean;
        description "If the neighbor on this link offers ZTP.";
    }
    leaf flood-port {
        type inet:port-number;
        description "The flooding port.";
    }
    leaf lie-rx-port {
        type inet:port-number;
        description "The port of LIE packet receiving.";
    }
    leaf holdtime {
        type rt-types:timer-value-seconds16;
        units seconds;
        description "The holding time of this adjacency.";
    }
    leaf local-nonce {
        type uint16;
        description "Local Nonce of the adjacency as advertised in LIEs.
                    In case of LIE packet this MUST correspond to the
                    value in the serialized object otherwise the packet
                    MUST be discarded.";
    }

    description "The interface information on this node.";
}
description "The interface information.";
}

grouping prefix-info {
    leaf prefix {
        type inet:ip-prefix;
        description "The prefix information.";
    }
}
```

```
    }
    leaf metric {
        type uint32;
        description "The metric of this prefix.";
    }
    leaf tag {
        type uint64;
        description "The tag of this prefix.";
    }
    leaf monotonic_clock {
        type PrefixSequenceType;
        description "The monotonic clock for mobile addresses.";
    }
    leaf from-link {
        type linkidtype;
        description "In case of locally originated prefixes, i.e.
            interface addresses this can describe which
            link the address belongs to.";
    }

    description "The detail information of prefix.";
}

grouping tie-id {
    leaf originator {
        type systemid;
        description "The originator's systemid of this TIE.";
    }
    uses direction-type;
    leaf tie-number {
        type uint32;
        description "The number of this TIE";
    }
    description "TIE is the acronym for 'Topology Information Element'.
        TIEs are exchanged between RIFT nodes to describe parts
        of a network such as links and address prefixes. This is
        the TIE identification information.";
}

grouping key-value {
    leaf key {
        type uint16;
        description "The type of key value combination.";
    }
    leaf value {
        type uint32;
        description "The value of key value combination.";
    }
}
```

```

    description "The key-value store information.";
  }

  grouping tie-info {
    leaf seq {
      type uint32;
      description "The sequence number of a TIE.";
    }
    leaf lifetime {
      type uint16 {
        range "1 .. 65535";
      }
      description "The lifetime of a TIE.";
    }
  }

  container tie-node {
    uses base-node-info;
    description "The node element information in this TIE.";
  }

  container tie-prefix {
    container prefixes {
      uses prefix-info;
      description "It is the prefixes TIE element.";
    }
    container positive_disaggregation_prefixes {
      uses prefix-info;
      description "It is the positive disaggregation prefixes TIE element.";
    }
    container negative_disaggregation_prefixes {
      uses prefix-info;
      description "It is the negative disaggregation prefixes element.";
    }
    container external_prefixes {
      uses prefix-info;
      description "It is the external prefixes element.";
    }
    description "The prefix information in this TIE.";
  }

  container kvs {
    uses key-value;
    description "The key/values in the database.";
  }

  description "TIE is the acronym for 'Topology Information Element'.
    TIEs are exchanged between RIFT nodes to describe parts
    of a network such as links and address prefixes. This TIE
    info is used to indicate the state of this TIE. When the
    type of this TIE is set to 'node', the node-element is

```



```

        making sense. When the type of this TIE is set to other
        types except for 'node', the prefix-info is making sense.";
    }

    grouping algorithm {
        choice algorighm-type {
            case spf {
                description "The algorithm is SPF.";
            }
            case all-path {
                description "The algorithm is all-path.";
            }
            description "The possible algorithm types.";
        }
        description "The computation algorithm types.";
    }

    /*
    * Data nodes
    */
    augment "/rt:routing/rt:control-plane-protocols/rt:control-plane-protocol" {
        when "derived-from-or-self(rt:type, 'rift:rift')" {
            description "This augment is only valid for a routing protocol instance of RIFT.";
        }
        description "RIFT ( Routing in Fat Trees ) YANG model.";
        container rift {
            presence "Container for RIFT protocol.";
            description "RIFT configuration data.";

            container node-info {
                description "The node information about RIFT.";
                uses node;
            }
            container neighbor {
                config false;
                list nbrs {
                    key "systemid remote-id";
                    uses base-node-info;
                    uses neighbor;
                    description "The information of a neighbor.";
                }
                description "The neighbor's information.";
            }
        }
        container database {
            config false;
            list ties {
                key "tie-index";
            }
        }
    }

```

```

        leaf tie-index {
            type uint32;
            description "The index of a TIE.";
        }
        container database-tie {
            uses tie-id;
            uses tie-info;

            description "The TIEs in the database.";
        }
        description "The detail information of a TIE.";
    }
    description "The TIEs information in database.";
} //database

container kv-store {
    list kvs {
        key "kvs-index";
        leaf kvs-index {
            type uint32;
            description "The index of a kv pair.";
        }

        container kvs-tie {
            uses tie-id;
            uses key-value;
            description "The TIEs in the kv-store.";
        }
        description "The information used to distinguish a Key/Value
pair.
ent is
er values
When the type of kv is set to 'node', node-elm
making sense. When the type of kv is set to oth
except 'node', prefix-info is making sense.";
    }
    description "The Key/Value store information.";
} //kv-store

} //rift
} //augment

/*
 * RPCs
 */

/*
 * Notifications
 */

```

```
notification error-set {
  description "The errors notification of RIFT.";
  container tie-level-error {
    uses tie-id;
    uses tie-info;
    description "The level is undefined in the LIEs.";
  }
  container nbr-error {
    list nbrs {
      key "systemid remote-id";
      uses base-node-info;
      uses neighbor;
      description "The information of a neighbor.";
    }
    description "The neighbor errors set.";
  }
}
}
<CODE ENDS>
```

## 8. Security Considerations

The YANG module specified in this document defines a schema for data that is designed to be accessed via network management protocols such as NETCONF [RFC6241] or RESTCONF [RFC8040]. The lowest NETCONF layer is the secure transport layer, and the mandatory-to-implement secure transport is Secure Shell (SSH) [RFC6242]. The lowest RESTCONF layer is HTTPS, and the mandatory-to-implement secure transport is TLS [RFC5246].

The NETCONF access control model [RFC6536] provides the means to restrict access for particular NETCONF or RESTCONF users to a preconfigured subset of all available NETCONF or RESTCONF protocol operations and content.

There are a number of data nodes defined in this YANG module that are writable/creatable/deletable (i.e., config true, which is the default). These data nodes may be considered sensitive or vulnerable in some network environments. Write operations (e.g., edit-config) to these data nodes without proper protection can have a negative effect on network operations.

The RPC operations in this YANG module may be considered sensitive or vulnerable in some network environments. It is thus important to control access to these operations.

## 9. IANA Considerations

The IANA is requested to assign two new URIs from the IETF XML registry ([RFC3688]). Authors are suggesting the following URI:

URI: urn:ietf:params:xml:ns:yang:ietf-rift

Registrant Contact: RIFT WG

XML: N/A, the requested URI is an XML namespace

This document also requests one new YANG module name in the YANG Module Names registry ([RFC6020]) with the following suggestion:

name: ietf-rift

namespace: urn:ietf:params:xml:ns:yang:ietf-rift

prefix: rift

reference: RFC XXXX

## 10. Contributors

The authors would like to thank Tony Przygienda, Benchong Xu (xu.benchong@zte.com.cn), for their review and valuable contributions.

## 11. Normative References

- [I-D.ietf-rift-rift] Team, T., "RIFT: Routing in Fat Trees", draft-ietf-rift-rift-05 (work in progress), April 2019.
- [I-D.ietf-rtgwg-policy-model] Qu, Y., Tantsura, J., Lindem, A., and X. Liu, "A YANG Data Model for Routing Policy Management", draft-ietf-rtgwg-policy-model-06 (work in progress), March 2019.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC5246] Dierks, T. and E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.2", RFC 5246, DOI 10.17487/RFC5246, August 2008, <<https://www.rfc-editor.org/info/rfc5246>>.

- [RFC6020] Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.
- [RFC6087] Bierman, A., "Guidelines for Authors and Reviewers of YANG Data Model Documents", RFC 6087, DOI 10.17487/RFC6087, January 2011, <<https://www.rfc-editor.org/info/rfc6087>>.
- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6242] Wasserman, M., "Using the NETCONF Protocol over Secure Shell (SSH)", RFC 6242, DOI 10.17487/RFC6242, June 2011, <<https://www.rfc-editor.org/info/rfc6242>>.
- [RFC6536] Bierman, A. and M. Bjorklund, "Network Configuration Protocol (NETCONF) Access Control Model", RFC 6536, DOI 10.17487/RFC6536, March 2012, <<https://www.rfc-editor.org/info/rfc6536>>.
- [RFC6991] Schoenwaelder, J., Ed., "Common YANG Data Types", RFC 6991, DOI 10.17487/RFC6991, July 2013, <<https://www.rfc-editor.org/info/rfc6991>>.
- [RFC7223] Bjorklund, M., "A YANG Data Model for Interface Management", RFC 7223, DOI 10.17487/RFC7223, May 2014, <<https://www.rfc-editor.org/info/rfc7223>>.
- [RFC7277] Bjorklund, M., "A YANG Data Model for IP Management", RFC 7277, DOI 10.17487/RFC7277, June 2014, <<https://www.rfc-editor.org/info/rfc7277>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC8040] Bierman, A., Bjorklund, M., and K. Watsen, "RESTCONF Protocol", RFC 8040, DOI 10.17487/RFC8040, January 2017, <<https://www.rfc-editor.org/info/rfc8040>>.
- [RFC8177] Lindem, A., Ed., Qu, Y., Yeung, D., Chen, I., and J. Zhang, "YANG Data Model for Key Chains", RFC 8177, DOI 10.17487/RFC8177, June 2017, <<https://www.rfc-editor.org/info/rfc8177>>.

- [RFC8342] Bjorklund, M., Schoenwaelder, J., Shafer, P., Watsen, K., and R. Wilton, "Network Management Datastore Architecture (NMDA)", RFC 8342, DOI 10.17487/RFC8342, March 2018, <<https://www.rfc-editor.org/info/rfc8342>>.
- [RFC8349] Lhotka, L., Lindem, A., and Y. Qu, "A YANG Data Model for Routing Management (NMDA Version)", RFC 8349, DOI 10.17487/RFC8349, March 2018, <<https://www.rfc-editor.org/info/rfc8349>>.
- [RFC8407] Bierman, A., "Guidelines for Authors and Reviewers of Documents Containing YANG Data Models", BCP 216, RFC 8407, DOI 10.17487/RFC8407, October 2018, <<https://www.rfc-editor.org/info/rfc8407>>.

## Authors' Addresses

Zheng Zhang  
ZTE Corporation

Email: [zhang\\_z@zte.com.cn](mailto:zhang_z@zte.com.cn)

Yuehua Wei  
ZTE Corporation

Email: [wei.yuehua@zte.com.cn](mailto:wei.yuehua@zte.com.cn)

Shaowen Ma  
Mellanox

Email: [mashaowen@gmail.com](mailto:mashaowen@gmail.com)

Xufeng Liu  
Volta Networks

Email: [xufeng.liu.ietf@gmail.com](mailto:xufeng.liu.ietf@gmail.com)