

RIFT
Internet-Draft
Intended status: Standards Track
Expires: 10 January 2022

J. Head, Ed.
T. Przygienda
W. Lin
Juniper Networks
9 July 2021

RIFT Auto-EVPN
draft-head-rift-auto-evpn-01

Abstract

This document specifies procedures that allow an EVPN overlay to be fully and automatically provisioned when using RIFT as underlay by leveraging RIFT's no-touch ZTP architecture.

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 10 January 2022.

Copyright Notice

Copyright (c) 2021 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 | 3 |
| 1.1. Requirements Language | 3 |
| 2. Design Considerations | 3 |
| 3. System ID | 4 |
| 4. Fabric ID | 4 |
| 5. Auto-EVPN Device Roles | 5 |
| 5.1. All Participating Nodes | 5 |
| 5.2. ToF Nodes as Route Reflectors | 5 |
| 5.3. Leaf Nodes | 6 |
| 6. Auto-EVPN Variable Derivation | 7 |
| 6.1. Auto-EVPN Version | 8 |
| 6.2. MAC-VRF ID | 8 |
| 6.3. Loopback Address | 8 |
| 6.3.1. Leaf Nodes as Gateways | 8 |
| 6.3.2. ToF Nodes as Route Reflectors | 9 |
| 6.3.2.1. Route Reflector Election Procedures | 9 |
| 6.4. Autonomous System Number | 10 |
| 6.5. Router ID | 10 |
| 6.6. Cluster ID | 10 |
| 6.7. Route Target | 10 |
| 6.8. Route Distinguisher | 10 |
| 6.9. EVPN MAC-VRF Services | 11 |
| 6.9.1. Untagged Traffic in Multiple Fabrics | 11 |
| 6.9.1.1. VLAN | 11 |
| 6.9.1.2. VNI | 11 |
| 6.9.1.3. MAC Address | 11 |
| 6.9.1.4. IPv6 IRB Gateway Address | 12 |
| 6.9.1.5. IPv4 IRB Gateway Address | 12 |
| 6.9.2. Tagged Traffic in Multiple Fabrics | 12 |
| 6.9.2.1. VLAN | 12 |
| 6.9.2.2. VNI | 12 |
| 6.9.2.3. MAC Address | 13 |
| 6.9.2.4. IPv6 IRB Gateway Address | 13 |
| 6.9.2.5. IPv4 IRB Gateway Address | 13 |
| 6.9.3. Tagged Traffic in a Single Fabric | 13 |
| 6.9.3.1. VLAN | 13 |
| 6.9.3.2. VNI | 14 |
| 6.9.3.3. MAC Address | 14 |
| 6.9.3.4. IPv6 IRB Gateway Address | 14 |
| 6.9.3.5. IPv4 IRB Gateway Address | 14 |
| 6.9.4. Traffic Routed to External Destinations | 14 |
| 6.9.4.1. Route Distinguisher | 15 |
| 6.9.4.2. Route Target | 15 |
| 6.10. Auto-EVPN Analytics | 15 |
| 6.10.1. Auto-EVPN Global Analytics Key Type | 16 |
| 6.10.2. Auto-EVPN MAC-VRF Key Type | 16 |

| | |
|---|----|
| 7. Acknowledgements | 18 |
| 8. Security Considerations | 18 |
| 9. References | 18 |
| 9.1. Normative References | 18 |
| Appendix A. Thrift Models | 19 |
| A.1. RIFT LIE Schema | 19 |
| A.1.1. Auto-EVPN Version | 19 |
| A.1.2. Fabric ID | 19 |
| A.2. RIFT Node-TIE Schema | 19 |
| A.2.1. Auto-EVPN Version | 19 |
| A.2.2. Fabric ID | 19 |
| A.3. common_evpn.thrift | 19 |
| A.4. auto_evpn_kv.thrift | 22 |
| Appendix B. Auto-EVPN Variable Derivation | 24 |
| Authors' Addresses | 36 |

1. Introduction

RIFT is a protocol that focuses heavily on operational simplicity. [RIFT] natively supports Zero Touch Provisioning (ZTP) functionality that allows each node in an underlay network to automatically derive its place in the topology and configure itself accordingly when properly cabled. RIFT can also disseminate Key-Value information contained in Key-Value Topology Information Elements (KV-TIEs) [RIFT-KV]. These KV-TIEs can contain any information and therefore be used for any purpose. Leveraging RIFT to provision EVPN overlays without any need for configuration and leveraging KV capabilities to easily validate correct operation of such overlay without a single point of failure would provide significant benefit to operators in terms of simplicity and robustness of such a solution.

1.1. 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 2119 [RFC2119].

2. Design Considerations

EVPN supports various service models, this document defines a method for the VLAN-Aware service model defined in [RFC7432]. Other service models may be considered in future revisions of this document.

Each model has its own set of requirements for deployment. For example, a functional BGP overlay is necessary to exchange EVPN NLRI regardless of the service model. Furthermore, the requirements are made up of individual variables, such as each node's loopback address and AS number for the BGP session. Some of these variables may be

coordinated across each node in a network, but are ultimately locally significant (e.g. route distinguishers). Similarly, calculation of some variables will be local only to each device. RIFT contains currently enough topology information in each node to calculate all those necessary variables automatically.

Once the EVPN overlay is configured and becomes operational, RIFT Key-Value TIEs can be used to distribute state information to allow for validation of basic operational correctness without the need for further tooling.

3. System ID

The 64-bit RIFT System ID that uniquely identifies a node as defined in RIFT [RIFT].

4. Fabric ID

RIFT operates on variants of Clos substrate which are commonly called an IP Fabric. Since EVPN VLANs can be either contained within one fabric or span them, Auto-EVPN introduces the concept of a Fabric ID into RIFT.

This section describes an optional extension to LIE packet schema in the form of a 16-bit Fabric ID that identifies a nodes membership within a particular fabric. Auto-EVPN capable nodes MUST support this extension but MAY not advertise it when not participating in Auto-EVPN. A non-present Fabric ID and value of 0 is reserved as ANY_FABRIC and MUST NOT be used for any other purpose.

Fabric ID MUST be considered in existing adjacency FSM rules so nodes that support Auto-EVPN can interoperate with nodes that do not. The LIE validation is extended with following clause and if it is not met, miscabbling should be declared:

```
(if fabric_id is not advertised by either node OR
 if fabric_id is identical on both nodes)
  AND
 (if auto_evpn_version is not advertised by either node OR
  if auto_evpn_version is identical on both nodes)
```

The appendix details LIE (Appendix A.1.2) and Node-TIE (Appendix A.2.2) schema changes.

5. Auto-EVPN Device Roles

Auto-EVPN requires that each node understand its given role within the scope of the EVPN implementation so each node derives the necessary variables and provides the necessary overlay configuration. For example, a leaf node performing VXLAN gateway functions does not need to derive its own Cluster ID or learn one from the route reflector that it peers with.

5.1. All Participating Nodes

Not all nodes have to participate in Auto-EVPN, however if a node does assume an Auto-EVPN role, it MUST derive the following variables:

IPv6 Loopback Address

Unique IPv6 loopback address used in BGP sessions.

Router ID

The BGP Router ID.

Autonomous System Number

The ASN for IBGP sessions.

Cluster ID

The Cluster ID for Top-of-Fabric IBGP route reflection.

5.2. ToF Nodes as Route Reflectors

This section defines an Auto-EVPN role whereby some Top-of-Fabric nodes act as EVPN route reflectors. It is expected that route reflectors would establish IBGP sessions with leaf nodes in the same fabric. The typical route reflector requirements do not change, however determining which specific values to use requires further consideration. ToF nodes performing route reflector functionality MUST derive the following variables:

IPv6 RR Loopback Address

The source address for IBGP sessions with leaf nodes in case ToF won election for one of the route reflectors in the fabric.

IPv6 RR Acceptable Prefix Range

Range of addresses acceptable by the route reflector to form a IBGP session. This range covers ALL possible IPv6 Loopback Addresses derived by other Auto EVPN nodes in the current fabric and other Auto-EVPN RRs addresses.

5.3. Leaf Nodes

Leaf nodes derive their role from realizing they are at the bottom of the fabric, i.e. not having any southbound adjacencies. Alternately, a node can assume a leaf node if it has only southbound adjacencies to nodes with explicit LEAF_LEVEL to allow for scenarios where RIFT leaves do NOT participate in Auto-EVPN.

Leaf nodes MUST derive the following variables:

IPv6 RR Loopback Addresses

Addresses of the RRs present in the fabric. Those addresses are used to build BGP sessions to the RR.

EVis

Leaf node derives all the necessary variables to instantiate EVIs with layer-2 and optionally layer-3 functionality.

If a leaf node is required to perform layer-2 VXLAN gateway functions, it MUST be capable of deriving the following types of variables:

Route Distinguisher

The route distinguisher corresponding to a MAC-VRF that uniquely identifies each node.

Route Target

The route target that corresponds to a MAC-VRF.

MAC VRF Name

This is an optional variable to provide a common MAC VRF name across all leaves.

Set of VLANs

Those are VLANs provisioned either within the fabric or allowing to stretch across fabrics.

For each VLAN derived in an EVI the following variables MUST be derived:

VLAN

The VLAN ID.

Name

This is an optional variable to provide a common VLAN name across all leaves.

VNI

The VNI that corresponds to the VLAN ID. This will contribute to the EVPN Type-2 route.

IRB

Optional variables of the IRB for the VLAN if the leaf performs layer-3 gateway function.

If a leaf node is required to perform layer-3 VXLAN gateway functions, it MUST additionally be capable of deriving the following types of variables:

IP Gateway MAC Address

The MAC address associated with IP gateway.

IP Gateway Subnetted Address

The IPv4 and/or IPv6 gateway address including its subnet length.

Type-5 EVPN IP Prefix with ToFs performing gateway functionality can also be derived and will be described in a future version of this document.

6. Auto-EVPN Variable Derivation

As previously mentioned, not all nodes are required to derive all variables in a given network (e.g. a transit spine node may not need to derive any or participate in Auto-EVPN). Additionally, all derived variables are derived from RIFT's FSM or ZTP mechanism so no additional flooding beside RIFT flooding is necessary for the functionality.

It is also important to mention that all variable derivation is in some way based on combinations of System ID, MAC-VRF ID, Fabric ID, EVI and VLAN and MUST comply precisely with calculation methods specified in the Auto-EVPN Variable Derivation section to allow interoperability between different implementations. All foundational code elements such as imports, constants, etc. are also mentioned there.

6.1. Auto-EVPN Version

This section describes extensions to both the RIFT LIE packet and Node-TIE schemas in the form of a 16-bit value that identifies the Auto-EVPN Version. Auto-EVPN capable nodes MUST support this extension, but MAY choose not to advertise it in LIEs and Node-TIEs when Auto-EVPN is not being utilized. The appendix describes LIE (Appendix A.1.1) and Node-TIE (Appendix A.2.1) schema changes in detail.

6.2. MAC-VRF ID

This section describes a variable MAC-VRF ID that uniquely identifies an instance of EVPN instance (EVI) and is used in variable derivation procedures. Each EVPN EVI MUST be associated with a unique MAC-VRF ID, this document does not specify a method for making that association or ensuring that they are coordinated properly across fabric(s).

6.3. Loopback Address

First and foremost, RIFT does not advertise anything more specific than the fabric default route in the southbound direction by default. However, Auto-EVPN nodes MUST advertise specific loopback addresses southbound to all other Auto-EVPN nodes so to establish MP-BGP reachability correctly in all scenarios.

Auto-EVPN nodes MUST derive a ULA-scoped IPv6 loopback address to be used as both the IBGP source address, as well as the VTEP source when VXLAN gateways are required. Calculation is done using the 6-bytes of reserved ULA space, the 2-byte Fabric ID, and the node's 8-byte System ID. Derivation of the System ID varies slightly depending upon the node's location/role in the fabric and will be described in subsequent sections.

6.3.1. Leaf Nodes as Gateways

Calculation is done using the 6-bytes of reserved ULA space, the 2-byte Fabric ID, and the node's 8-byte System ID.

In order for leaf nodes to derive IPv6 loopback addresses, algorithms shown in both `auto_evpn_fidsidv6loopback` (Figure 24) and `auto_evpn_v6prefixfidsid2loopback` (Figure 9) are required.

IPv4 addresses MAY be supported, but it should be noted that they have a higher likelihood of collision. The appendix contains the required `auto_evpn_fidsid2v4loopback` (Figure 23) algorithm to support IPv4 loopback derivation.

6.3.2. ToF Nodes as Route Reflectors

ToF nodes acting as route reflectors MUST derive their loopback address according to the specific section describing the algorithm. Calculation is done using the 6-bytes of reserved ULA space, the 2-byte Fabric ID, and the 8-byte System ID of each elected route reflector.

In order for the ToF nodes to derive IPv6 loopbacks, the algorithms shown in both `auto_evpn_fidsidv6loopback` (Figure 24) and `auto_evpn_fidrrpref2rrloopback` (Figure 10) are required.

In order for the ToF derive the necessary prefix range to facilitate peering requests from any leaf, the algorithm shown in "`auto_evpn_fid2fabric_prefixes`" (Figure 8) is required.

6.3.2.1. Route Reflector Election Procedures

Four Top-of-Fabric nodes MUST be elected as an IBGP route reflector. Each ToF performs the election independently based on system IDs of other ToFs in the fabric obtained via southbound reflection. The route reflector election procedures are defined as follows:

1. ToF node with the highest System ID.
2. ToF node with the lowest System ID.
3. ToF node with the 2nd highest System ID.
4. ToF node with the 2nd lowest System ID.

This ordering is necessary to prevent a single node with either the highest or lowest System ID from triggering changes to route reflector loopback addresses as it would result in all BGP sessions dropping.

For example, if two nodes, ToF01 and ToF02 with System IDs 002c6af5a281c000 and 002c6bf5788fc000 respectively, ToF02 would be elected due to it having the highest System ID of the ToFs (002c6bf5788fc000). If a ToF determines that it is elected as route reflector, it uses the knowledge of its position in the list to derive route reflector v6 loopback address.

The algorithm shown in "`auto_evpn_sids2rrs`" (Figure 6) is required to accomplish this.

Considerations for multiplane route reflector elections will be included in future revisions.

6.4. Autonomous System Number

Nodes in each fabric MUST derive a private autonomous system number based on its Fabric ID so that it is unique across the fabric.

The algorithm shown in `auto_evpn_fid2private_AS` (Figure 25) is required to derive the private ASN.

6.5. Router ID

Nodes MUST derive a Router ID that is based on both its System ID and Fabric ID so that it is unique to both.

The algorithm shown in `auto_evpn_sidfid2bgpid` (Figure 11) is required to derive the BGP Router ID.

6.6. Cluster ID

Route reflector nodes in each fabric MUST derive a cluster ID that is based on its Fabric ID so that it is unique across the fabric.

The algorithm shown in `auto_evpn_fid2clusterid` (Figure 26) is required to derive the BGP Cluster ID.

6.7. Route Target

Nodes hosting EVPN EVIs MUST derive a route target extended community based on the MAC-VRF ID for each EVI so that it is unique across the network. Route targets MUST be of type 0 as per RFC4360.

For example, if given a MAC-VRF ID of 1, the derived route target would be "target:1"

The algorithm shown in `auto_evpn_evi2rt` (Figure 12) is required to derive the Route Target community.

6.8. Route Distinguisher

Nodes hosting EVPN EVIs MUST derive a type-0 route distinguisher based on its System ID and Fabric ID so that it is unique per MAC-VRF and per node.

The algorithm shown in `auto_evpn_sidfid2rd` (Figure 18) is required to derive the Route Distinguisher.

6.9. EVPN MAC-VRF Services

It's obvious that applications utilizing Auto-EVPN overlay services may require a variety of layer-2 and/or layer-3 traffic considerations. Variables supporting these services are also derived based on some combination of MAC-VRF ID, Fabric ID, and other constant values. Integrated Routing and Bridging (IRB) gateway address derivation also leverages a set of constant RANDOMSEEDS (Figure 5) values that MUST be used to provide additional entropy.

In order to ensure that VLAN ID's don't collide, a single deployment SHOULD NOT exceed 3 fabrics with 3 EVIs where each EVI terminate 15 VLANs. The algorithms shown in `auto_evpn_fidevivlansvlans2desc` (Figure 16) and `auto_evpn_vlan_description_table` (Figure 15) are required to derive VLANs accordingly. An implementation MAY exceed this, but MUST indicate methods to ensure collision-free derivation and describe which VLANs are stretched across fabrics.

6.9.1. Untagged Traffic in Multiple Fabrics

This section defines methods to derive unique VLAN, VNI, MAC, and gateway address values for deployments where untagged traffic is stretched across multiple fabrics.

6.9.1.1. VLAN

Untagged traffic stretched across multiple fabrics MUST derive VLAN tags based on MAC-VRF ID in conjunction with a constant value of 1 (i.e. MAC-VRF ID + 1).

6.9.1.2. VNI

Untagged traffic stretched across multiple fabrics MUST derive VNIs based on MAC-VRF ID and Fabric ID in conjunction with a constant value. These VNIs MUST correspond to EVPN Type-2 routes.

The algorithm shown in `auto_evpn_fidevivid2vni` (Figure 14) is required to derive VNIs for Type-2 EVPN routes.

6.9.1.3. MAC Address

The MAC address MUST be a unicast address and also MUST be identical for any IRB gateways that belong to an individual bridge-domain across fabrics. The last 5-bytes MUST be a hash of the MAC-VRF ID and a constant value of 1 that is calculated using the previously mentioned random seed values.

The algorithm shown in `auto_evpn_fidevivid2mac` (Figure 22) is required to derive MAC addresses.

6.9.1.4. IPv6 IRB Gateway Address

The derived IPv6 gateway address MUST be from a ULA-scoped range that will account for the first 6-bytes. The next 5-bytes MUST be the last bytes of the derived MAC address. Finally, the remaining 7-bytes MUST be `::0001`.

The algorithm shown in `auto_evpn_fidevivid2v6subnet` (Figure 21) is required to derive the IPv6 gateway address.

6.9.1.5. IPv4 IRB Gateway Address

The derived IPv4 gateway address MUST be from a RFC1918 range, which accounts for the first octet. The next octet MUST be a hash of the MAC-VRF ID and a constant value of 1 that is calculated using the previously mentioned random seed values. Finally, the remaining 2 octets MUST be 0 and 1 respectively.

The algorithm shown in `auto_evpn_v4prefixfidevivid2v4subnet` (Figure 19) is required to derive the IPv4 gateway address. It should be noted that there is a higher likelihood of address collisions when deriving IPv4 addresses.

6.9.2. Tagged Traffic in Multiple Fabrics

This section defines methods to derive unique VLAN, VNI, MAC, and gateway address values for deployments where tagged traffic is stretched across multiple fabrics.

6.9.2.1. VLAN

Tagged traffic stretched across multiple fabrics MUST derive VLAN tags based on MAC-VRF ID in conjunction with a constant value of 16 (i.e. `MAC-VRF ID + 16`).

6.9.2.2. VNI

Tagged traffic stretched across multiple fabrics MUST derive VNIs based on MAC-VRF ID and Fabric ID in conjunction with a constant value. These VNIs MUST correspond to EVPN Type-2 routes.

The algorithm shown in `auto_evpn_fidevivid2vni` (Figure 14) is required to derive VNIs for Type-2 EVPN routes.

6.9.2.3. MAC Address

The MAC address MUST be a unicast address and also MUST be identical for any IRB gateways that belong to an individual bridge-domain across fabrics. The last 5-bytes MUST be a hash of the MAC-VRF ID and a constant value of 1 that is calculated using the previously mentioned random seed values.

The algorithm shown in `auto_evpn_fidevividssid2mac` (Figure 22) is required to derive MAC addresses.

6.9.2.4. IPv6 IRB Gateway Address

The derived IPv6 gateway address MUST be from a ULA-scoped range that will account for the first 6-bytes. The next 5-bytes MUST be the last bytes of the derived MAC address. Finally, the remaining 7-bytes MUST be `::0001`.

The algorithm shown in `auto_evpn_fidevividssid2v6subnet` (Figure 21) is required to derive the IPv6 gateway address.

6.9.2.5. IPv4 IRB Gateway Address

The derived IPv4 gateway address MUST be from a RFC1918 range, which accounts for the first octet. The next octet MUST be a hash of the MAC-VRF ID and a constant value of 16 that is calculated using the previously mentioned random seed values. Finally, the remaining 2 octets MUST be 0 and 1 respectively.

The algorithm shown in `auto_evpn_v4prefixfidevividssid2v4subnet` (Figure 19) is required to derive the IPv4 gateway address. It should be noted that there is a higher likelihood of address collisions when deriving IPv4 addresses.

6.9.3. Tagged Traffic in a Single Fabric

This section defines a method to derive unique VLAN, VNI, MAC, and gateway address values for deployments where untagged traffic is contained within a single fabric.

6.9.3.1. VLAN

Tagged traffic contained to a single fabric MUST derive VLAN tags based on MAC-VRF ID and Fabric ID in conjunction with a constant value of 17 (i.e. MAC-VRF ID + Fabric ID + 17).

6.9.3.2. VNI

Tagged traffic contained to a single fabric MUST derive VNIs based on MAC-VRF ID and Fabric ID in conjunction with a constant value. These VNIs MUST correspond to EVPN Type-2 routes.

The algorithm shown in `auto_evpn_fidevivid2vni` (Figure 14) is required to derive VNIs for Type-2 EVPN routes.

6.9.3.3. MAC Address

The MAC address MUST be a unicast address and also MUST be identical for any IRB gateways that belong to an individual bridge-domain across fabrics. The last 5-bytes MUST be a hash of the MAC-VRF ID and a constant value of 1 that is calculated using the previously mentioned random seed values.

The algorithm shown in `auto_evpn_fidevividsid2mac` (Figure 22) is required to derive MAC addresses.

6.9.3.4. IPv6 IRB Gateway Address

The derived IPv6 gateway address MUST be from a ULA-scoped range, which accounts for the first 6-bytes. The next 5-bytes MUST be the last bytes of the derived MAC address. Finally, the remaining 7-bytes MUST be `::0001`.

The algorithm shown in `auto_evpn_fidevividsid2v6subnet` (Figure 21) is required to derive the IPv6 gateway address.

6.9.3.5. IPv4 IRB Gateway Address

The derived IPv4 gateway address MUST be from a RFC1918 range, which accounts for the first octet. The next octet MUST be a hash of the MAC-VRF ID and a constant value of 17 that is calculated using the previously mentioned random seed values. Finally, the remaining 2 octets MUST be 0 and 1 respectively.

The algorithm shown in `auto_evpn_v4prefixfidevividsid2v4subnet` (Figure 19) is required to derive the IPv4 gateway address. It should be noted that there is a higher likelihood of address collisions when deriving IPv4 addresses.

6.9.4. Traffic Routed to External Destinations

6.9.4.1. Route Distinguisher

Nodes hosting IP Prefix routes MUST derive a type-0 route distinguisher based on its System ID and Fabric ID so that it is unique per IP-VRF and per node.

The algorithm shown in `auto_evpn_sidfid2rd` (Figure 18) is required to derive the Route Target.

6.9.4.2. Route Target

Nodes hosting IP prefix routes MUST derive a route target extended community based on the MAC-VRF ID for each IP-VRF so that it is unique across the network. Route targets MUST be of type 0.

The algorithm shown in `auto_evpn_evi2rt` (Figure 12) is required to derive the Route Target community.

6.10. Auto-EVPN Analytics

Leaf nodes MAY optionally advertise analytics information about the Auto-EVPN fabric to ToF nodes using RIFT Key-Value TIEs. This may be advantageous in that overlay validation and troubleshooting activities can be performed on the ToF nodes.

This section requests suggested values from the RIFT Well-Known Key-Type Registry and describes their use for Auto-EVPN.

| Name | Value | Description |
|-----------------------------|-------|--|
| Auto-EVPN Analytics MAC-VRF | 3 | Analytics describing a MAC-VRF on a particular node within a fabric. |
| Auto-EVPN Analytics Global | 4 | Analytics describing an Auto-EVPN node within a fabric. |

Table 1: Requested RIFT Key Registry Values

The normative Thrift schema can be found in the appendix (Appendix A.4).

6.10.1. Auto-EVPN Global Analytics Key Type

This Key Type describes node level information within the context of the Auto-EVPN fabric. The System ID of the advertising leaf node MUST be used to differentiate the node among other nodes in the fabric.

The Auto-EVPN Global Key Type MUST be advertised with the RIFT Fabric ID encoded into the 3rd and 4th bytes of the Key Identifier.

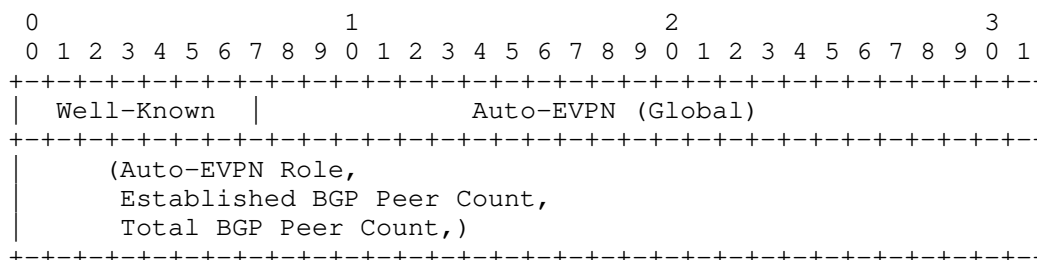


Figure 1: Auto-EVPN Global Key-Value TIE

where:

Auto-EVPN Role:

The value indicating the node's Auto-EVPN role within the fabric.

0: Illegal value, MUST NOT be used.

1: Auto-EVPN Leaf Gateway

2: Auto-EVPN Top-of-Fabric Gateway

Established BGP Session Count:

A 16-bit integer indicating the number of BGP sessions in the Established state.

Total BGP Peer Count:

A 16-bit integer indicating the total number of possible BGP sessions on the local node, regardless of state.

6.10.2. Auto-EVPN MAC-VRF Key Type

This Key-Value structure contains information about a specific MAC-VRF within the Auto-EVPN fabric.

The Auto-EVPN MAC-VRF Key Type MUST be advertised with the Auto-EVPN MAC-VRF ID encoded into the 3rd and 4th bytes of the Key Identifier.

All values advertised in a MAC-VRF Key-Value TIE MUST represent only state of the local node.

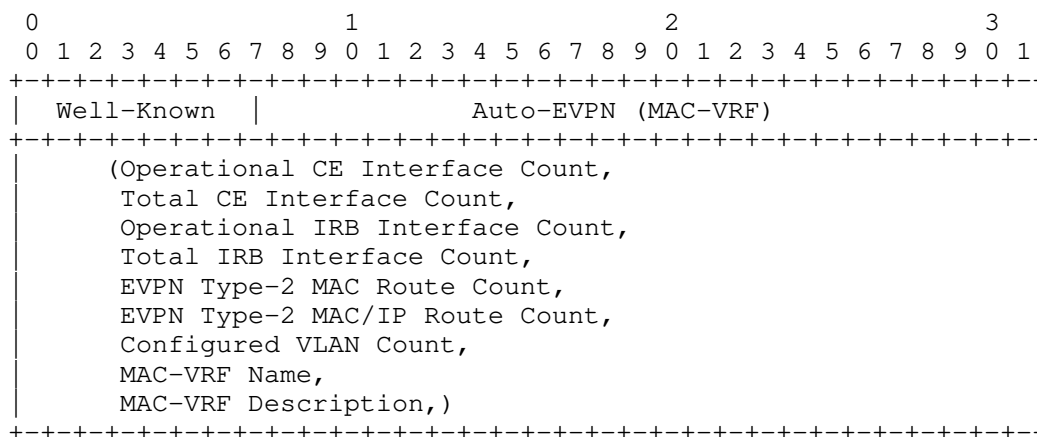


Figure 2: Auto-EVPN MAC-VRF Key-Value TIE

where:

Operational Customer Edge Interface Count:

A 16-bit integer indicating the number of CE interfaces associated with the MAC-VRF where both administrative and operational status are "up".

Total Customer Edge Interface Count:

A 16-bit integer indicating the total number of CE interfaces associated with the MAC-VRF regardless of interface status.

Operational IRB Interface Count:

A 16-bit integer indicating the number of IRB interfaces associated with the MAC-VRF where both administrative and operational status are "up".

Total IRB Interface Count:

A 16-bit integer indicating the total number of IRB interfaces associated with the MAC-VRF regardless of interface status.

EVPN Type-2 MAC Route Count:

A 32-bit integer indicating the total number of EVPN Type-2 MAC routes.

EVPN Type-2 MAC/IP Route Count:

A 32-bit integer indicating the total number of EVPN Type-2 MAC/IP routes.

VLAN Count:

A 16-bit integer indicating the total number configured VLANs.

MAC-VRF Name:

A string used to indicate the name of the MAC-VRF on the node.

MAC-VRF Description:

A string used to describe the MAC-VRF on the node, similar to that of an interface description.

7. Acknowledgements

The authors would like to thank Olivier Vandezande, Matthew Jones, and Michal Styszynski for their contributions.

8. Security Considerations

This document introduces no new security concerns to RIFT or other specifications referenced in this document.

9. References

9.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC7432] Sajassi, A., Aggarwal, R., Bitar, N., Isaac, A., Uttaro, J., Drake, J., and W. Henderickx, "BGP MPLS-Based Ethernet VPN", February 2015, <<https://www.rfc-editor.org/info/rfc7432>>.
- [RIFT] Przygienda, T., Sharma, A., Thubert, P., Rijsman, B., and D. Afanasiev, "RIFT: Routing in Fat Trees", Work in Progress, draft-ietf-rift-rift-13, July 2021.
- [RIFT-KV] Head, J. and T. Przygienda, "RIFT Keys Structure and Well-Known Registry in Key Value TIE", Work in Progress, draft-head-rift-kv-registry-01, July 2021.

Appendix A. Thrift Models

This section contains the normative Thrift models required to support Auto-EVPN. Per the main RIFT [RIFT] specification, all signed values MUST be interpreted as unsigned values.

A.1. RIFT LIE Schema

A.1.1. Auto-EVPN Version

```
struct LIEPacket {  
    ...  
    /** It provides optional version of EVPN ZTP as 256 * MAJOR + MINOR */  
    26: optional i16          auto_evpn_version;  
    ...
```

A.1.2. Fabric ID

```
struct LIEPacket {  
    ...  
    /** It provides the optional ID of the configured fabric */  
    25: optional common.FabricIDType fabric_id;  
    ...
```

A.2. RIFT Node-TIE Schema

A.2.1. Auto-EVPN Version

```
struct NodeTIEElement {  
    ...  
    /** It provides optional version of EVPN ZTP as 256 * MAJOR + MINOR */  
    13: optional i16          auto_evpn_version;  
    ...
```

A.2.2. Fabric ID

```
struct NodeTIEElement {  
    ...  
    /** It provides the optional ID of the Fabric configured */  
    12: optional common.FabricIDType fabric_id;  
    ...
```

A.3. common_evpn.thrift

This section contains the normative Auto-EVPN Thrift schema.

```
/**
    Thrift file for common AUTO EVPN definitions for RIFT

    Copyright (c) Juniper Networks, Inc., 2016-
    All rights reserved.
*/

namespace py common_evpn
namespace rs models

include "common.thrift"
include "encoding.thrift"
include "statistics.thrift"

const common.FabricIDType    default_fabric_id    = 1
const i32                    default_evis          = 3
const i32                    default_vlans_per_evi = 7

typedef i32      RouterIDType
typedef i32      ASType
typedef i32      ClusterIDType

struct EVPNAnyRole {
    1: required    common.IpV6Address              v6_loopback,
    2: required    common.IpV6Address              type5_v6_loopback,
    3: required    common.IpV4Address              type5_v4_loopback,
    4: required    RouterIDType                    bgp_router_id,
    5: required    ASType                          autonomous_system,
    6: required    ClusterIDType                   cluster_id,
    /** prefixes to be redistributed north */
    7: required    set<common.IPPrefixType>         redistribute_north,
    /** prefixes to be redistributed south */
    8: required    set<common.IPPrefixType>         redistribute_south,
    /** group name for evpn auto overlay */
    9: required    string                          bgp_group_name,
    /** fabric prefixes to be advertised in rift instead of default */
    10: required   set<common.IPPrefixType>         fabric_prefixes,
}

struct PartialEVPNEVI {
    // route target per RFC4360
    1: required    CommunityType                   rt_target,
    2: required    RTDistinguisherType             rt_distinguisher,
    3: required    RTDistinguisherType             rt_type5_distinguisher,
    5: required    string                          mac_vrf_name,
    6: required    VNIType                         type5_vni,
}
```

```

struct EVPNRRRole {
    2: required    common.Ipv6Address          v6_rr_addr_loopback,
    3: required    common.Ipv6PrefixType       v6_peers_allowed_range,
    4: required    map<MACVRFNumberType, PartialEVPNEVI> evis,
}

typedef i64      RTDistinguisherType
typedef i64      RTTargetType
typedef i16      MACVRFNumberType

typedef i16      VLANIDType
typedef binary   MACType

typedef i16      UnitType

struct IRBType {
    1: required    string                      name,
    2: required    UnitType                    unit,
    /// constant
    3: required    MACType                     mac,
    /// contains address of the gateway as well
    4: optional    common.Ipv6PrefixType       v6_subnet,
    /// contains address of the gateway as well
    5: optional    common.Ipv4PrefixType       v4_prefix,
}

typedef i32      VNIType

struct VLANType {
    1: optional    VLANIDType                  id,
    2: required    string                      name,
    3: optional    IRBType                     irb,
    5: optional    bool                        stretched = false,
    6: optional    bool                        is_native = false,
}

struct CEInterfaceType {
    2: optional    common.IEEE802_1ATimeStampType moved_to_ce,
    // we may not be able to obtain it in case of internal errors
    3: optional    string                      platform_interface_name,
}

typedef i64      CommunityType

struct EVPNEVI {
    // route target per RFC4360
    1: required    CommunityType               rt_target,
    2: required    RTDistinguisherType         rt_distinguisher,
}

```

```

        3: required    RTDistinguisherType          rt_type5_distinguisher,
        4: required    string                      mac_vrf_name,
        // fabric unique 24 bits VNI on non-stretch, otherwise unique across fabrics
        5: required    map<VNIType, VLANType>       vlans,
        6: required    VNIType                     type5_vni,
    }

    struct EVPNLeafRole {
        1: required    set<common.IPv6Address>      rrs,
        2: required    map<MACVRFNumberType, EVPNEVI> evis,
        3: optional    map<common.LinkIDType,
                        CEInterfaceType>            ce_interfaces,

        5: optional    binary                      leaf_unique_lacp_system_i
    d,
        6: optional    binary                      fabric_unique_lacp_system
    _id,
    }

    /// structure to indicate EVPN roles assumed and their variables for
    /// external platform to configure itself accordingly. Presence of
    /// according structure indicates that the role is assumed.
    struct EVPNRoles {
        1: required    EVPNAnyRole                  generic,
        2: optional    EVPNRRRole                   route_reflector,
        3: optional    EVPNLeafRole                 leaf,
    }

    const common.TimeIntervalInSecType    default_leaf_delay = 120
    const common.TimeIntervalInSecType    default_interface_ce_delay = 180
    /// default delay before EVPNZTP FSM starts to compute anything
    const common.TimeIntervalInSecType    default_evpnztp_startup_delay = 60

```

A.4. auto_evpn_kv.thrift

This section contains the normative Auto-EVPN Analytics Thrift schema.

```

include "common.thrift"

namespace py auto_evpn_kv
namespace rs models

/** We don't need the full role structure, only an indication of the node's basic
    role */
enum AutoEVPNRole {
    ILLEGAL          = 0,
    auto_evpn_leaf_erb = 1,
    auto_evpn_tof_gw   = 2,
}

```

```

enum    KVTypes {
    OUI          = 1,
    WellKnown    = 2,
}

const i8          AutoEVPNWellKnownKeyType  = 1
typedef i32        AutoEVPNKeyIdentifier
typedef i16        AutoEVPNCOUNTERType
typedef i32        AutoEVPNLongCounterType

const i8          GlobalAutoEVPNTelemetryKV = 4
const i8          AutoEVPNTelemetryKV      = 3

/** Per the according RIFT draft the key comes from the well known space.
    Part of the key is used as Fabric-ID.

    1st      byte  MUST be = "Well-Known"
    2nd      byte  MUST be = "Global Auto-EVPN Telemetry KV",
    3rd/4th bytes MUST be = FabricIDType
*/
struct AutoEVPNTelemetryGlobalKV {
    /** Only values that the ToF cannot derive itself should be flooded. */
    1: required    set<AutoEVPNRole>          auto_evpn_roles,

    /** Established BGP peer count (for Auto-EVPN)
    2: optional    AutoEVPNCOUNTERType        established_bgp_peer_count,

    /** Total BGP peer count (for Auto-EVPN)
    3: optional    AutoEVPNCOUNTERType        total_bgp_peer_count,
}

/** Per the according RIFT draft the key comes from the well known space.
    Part of the key is used as MAC-VRF number.

    1st      byte  MUST be = "Well-Known"
    2nd      byte  MUST be = indicates "Auto-EVPN Telemetry KV",
    3rd/4th bytes MUST be = MACVRFNumberType
*/
struct AutoEVPNTelemetryMACVRFKV {
    /** Active CE interface count (up/up)
    1: optional    AutoEVPNCOUNTERType        active_ce_interfaces,

    /** Total CE interface count
    2: optional    AutoEVPNCOUNTERType        total_ce_interfaces,

    /** Active IRB interface count (up/up)
    3: optional    AutoEVPNCOUNTERType        active_irb_interfaces,

```

```
/** Total IRB interface count
4: optional  AutoEVPNCounterType          total_irb_interfaces,

/** Local EVPN Type-2 MAC route count
5: optional  AutoEVPNLongCounterType       local_evpn_type2_mac_routes,

/** Local EVPN Type-2 MAC/IP route count
6: optional  AutoEVPNLongCounterType       local_evpn_type2_mac_ip_routes,

/** number of configured VLANs */
7: optional  il6                           configured_vlans,

/** optional human readable name */
8: optional  string                        name,

/** optional human readable string describing the MAC-VRF */
9: optional  string                        description,
}
```

Figure 3: Auto-EVPN Key-Value Thrift Schema

Appendix B. Auto-EVPN Variable Derivation


```

use std::cell::{RefCell, RefMut};
use std::cmp::{max, min};
use std::collections::{BTreeMap, BTreeSet, HashMap};
use std::fmt::Debug;
use std::net::{Ipv4Addr, Ipv6Addr};
use std::str::FromStr;
use itertools::interleave;
use itertools::Itertools;
use rayon::slice::ParallelSliceMut;

use foundation::models::common::{FabricIDType, IPv6PrefixType};
use foundation::models::common::LevelType;
use foundation::models::common::HierarchyIndications;
use foundation::models::common::IPPrefixType;
use foundation::models::common::IPv4Address;
use foundation::models::common::IPv4PrefixType;
use foundation::models::common::IPv6Address;
use foundation::models::common::LEAF_LEVEL;
use foundation::models::common_services::ServiceErrorType;
use foundation::models::common_evpn::{DEFAULT_EVIS, DEFAULT_VLANS_PER_EVI,
                                        DEFAULT_EVPNZTP_STARTUP_DELAY, DEFAULT_
FABRIC_ID, DEFAULT_INTERFACE_CE_DELAY,
                                        DEFAULT_LEAF_DELAY, EVPNAnyRole, EVPNLe
afRole,
                                        EVPNRoles, EVPNRRRole, UnitType};
use foundation::models::common_evpn::CEInterfaceType;
use foundation::models::common_evpn::CommunityType;
use foundation::models::common_evpn::EVPNEVI;
use foundation::models::common_evpn::IRBType;
use foundation::models::common_evpn::MACVRFNumberType;
use foundation::models::common_evpn::PartialEVPNEVI;
use foundation::models::common_evpn::RTDistinguisherType;
use foundation::models::common_evpn::VLANIDType;
use foundation::models::common_evpn::VLANType;
use foundation::models::common_evpn::VNIType;
use ILLEGAL_SYSTEM_ID;
use NodeCapabilities;
use UnsignedSystemID;

```

Figure 4: auto_evpn_imports

```

/// indicates how many RRs we're computing in EVPN ZTP
pub const MAX_AUTO_EVPN_RRS: usize = 3;
/// indicates the fabric has no ID, used in computations to omit effects of fabri
c ID
pub const NO_FABRIC_ID: FabricIDType = 0;
/// invalid MACVRF number, MACVRFs start from 1
pub const NO_MACVRF: MACVRFNumberType = 0;

/// unique v6 prefix for all nodes starts with this
pub const AUTO_EVPN_V6PREF: &str = "FD00:A1";

```

```
/// how many bytes in a v6pref for RRs
pub const AUTO_EVPN_V6PREFLEN: usize = 8 * 3;
/// unique v6 prefix for route reflector purposes starts like this
pub const AUTO_EVPN_V6RRPREF: &str = "FD00:A2";
/// unique v6 prefix for type-5 purposes starts like this
pub const AUTO_EVPN_V6T5PREF: &str = "FD00:A3";
/// unique v6 prefix for IRB prefix purposes
pub const AUTO_EVPN_V6IRBPREF: &str = "FD00";
/// unique v6 prefix first byte for v6 IRB prefix purposes
pub const AUTO_EVPN_V6IRBPREFFIRSTBYTE: u8 = 0xA4;
/// unique v4 prefix for IRB purposes
pub const AUTO_EVPN_V4IRBPREF: &str = "10";
/// 3 bytes of prefix type and then we have fabric ID after that
pub const AUTO_EVPN_V6_FABPREFIXLEN: usize = 8 + 8 + 8 + 16;

/// per RFC magic
const RT_TARGET_HIGH: CommunityType = 0;
const RT_TARGET_LOW: CommunityType = 0;

/// first available VLAN number
pub const FIRST_VLAN: VLANIDType = 1;
/// maximum vlan number one less than maximum to use as bitmask
pub const MAX_VLAN: VLANIDType = 4095;
/// constant VLAN shift
pub const FIRST_VLAN_SHIFT: VLANIDType = 16;
/// NATIVE VLAN number
pub const NATIVE_VLAN: VLANIDType = 1;

/// abstract description of VLAN properties for a derived VLAN
pub struct VLANDescription {
    pub vlan_id: VLANIDType,
    pub name: String,
    /// can this VLAN be stretched across multiple fabrics
    pub stretchable: bool,
    pub native: bool,
}

/// maximum number of VLANs per MACVRF
pub const MAX_VLANS_PER_EVI: usize = 15;

pub type VLANStretchableType = bool;
pub type VLANNativeType = bool;

pub const EXTRATYPE5_RD_DISTINGUISHER: u32 = 0xffff_ffff;

/// high bits of type 5 VNI
const TYPE5VNIHIGH: VNIDType = 0x0080_0000;
/// bitmask for type 2 VNI
```

```
const TYPE2VNIMASK: VNIType = 0x00ff_ffff ^ TYPE5VNIHIGH;

/// random seeds used in several algorithms to increase entropy
pub const RANDOMSEEDS: [u64; 4] = [
    27008318799u64,
    67438371571,
    37087353685,
    88675895388,
];
```

Figure 5: auto_evpn_const_structs_type

```
pub(crate) fn auto_evpn_sids2rrs(mut v: Vec<UnsignedSystemID>) -> Vec<UnsignedSystemID> {
    v.par_sort_unstable();
    let r = if v.len() > 2 {
        let mut s = v.split_off(v.len() / 2);
        s.reverse();
        interleave(v.into_iter(), s.into_iter()).collect()
    } else {
        v
    };
    r
}
```

Figure 6: auto_evpn_sids2rrs

```
pub(crate) fn auto_evpn_v62octets(a: Ipv6Addr) -> Vec<u8> {
    a.octets().iter().cloned().collect()
}
```

Figure 7: auto_evpn_v62octets

```

/// fabric prefixes derived instead of advertising default on the fabric to allow
/// for default route on ToF or leaves
pub fn auto_evpn_fid2fabric_prefixes(fid: FabricIDType) -> Result<Vec<IPPrefixType>
e>, ServiceErrorType> {
    vec![
        (auto_evpn_fidsidv6loopback(fid, ILLEGAL_SYSTEM_I_D as _), AUTO_EVPN_V6PR
EFLEN),
        (auto_evpn_fidrrpref2rrloopback(fid, ILLEGAL_SYSTEM_I_D as _), AUTO_EVPN_
V6PREFLEN),
    ]
    .into_iter()
    .map(|(p, _)|
        match p {
            Ok(_) => Ok(
                IPPrefixType::Ipv6prefix(
                    Ipv6PrefixType {
                        address: auto_evpn_v62octets(p?),
                        prefixlen: AUTO_EVPN_V6PREFLEN as _,
                    }),
            Err(e) => Err(e),
        }
    )
    .collect::()
}

```

Figure 8: auto_evpn_fid2fabric_prefixes

```

/// local address with encoded fabric ID and system ID for collision free identif
iers. Basis
/// for several different prefixes.
pub fn auto_evpn_v6prefixfidsid2loopback(v6pref: &str, fid: FabricIDType,
sid: UnsignedSystemID) -> Result<Ipv6Addr
r, ServiceErrorType> {
    assert!(fid != 0);
    let a = format!("{:02X}::{}",
        v6pref,
        fid as u16,
        sid.to_ne_bytes()
            .iter()
            .chunks(2)
            .into_iter()
            .map(|chunk|
                chunk.fold(0u16, |v, n| (v << 8) | *n as u16))
            .map(|v| format!("{:04X}", v))
            .collect::

```

Figure 9: auto_evpn_v6prefixfidsid2loopback

```

/// auto evpn V6 loopback for RRs
pub fn auto_evpn_fidrrpref2rrloopback(fid: FabricIDType,
                                     preference: u8) -> Result<Ipv6Addr, Service
ErrorType> {
    auto_evpn_v6prefixfidsid2loopback(AUTO_EVPN_V6RRPREF, fid, (1 + preference) as _)
}

```

Figure 10: auto_evpn_fidrrpref2rrloopback

```

/// auto evpn BGP router ID
pub fn auto_evpn_sidfid2bgpid(fid: FabricIDType, sid: UnsignedSystemID) -> u32 {
    assert!(fid != 0);
    let hs: u32 = ((sid & 0xffff_ffff_0000_0000) >> 32) as _;
    let mut ls: u32 = (sid & 0xffff_ffff) as _;
    ls = ls.rotate_right(7) ^ (fid as u32).rotate_right(13);
    max(1, hs ^ ls) // never a 0
}

```

Figure 11: auto_evpn_sidfid2bgpid

```

/// route target bytes are type0/0 and then add EVI
pub fn auto_evpn_evi2rt(evi: MACVRFNumberType) -> CommunityType {
    let wideevi = (evi + 1) as CommunityType;

    (RT_TARGET_HIGH << (64 - 8)) | (RT_TARGET_LOW << 64 - 16) |
    ((wideevi) << 17) |
    ((wideevi))
}

```

Figure 12: auto_evpn_evi2rt

```

/// type-5 VNI for an EVI
pub fn auto_evpn_fidevi2type5vni(fid: FabricIDType, evi: MACVRFNumberType) -> VNI
Type {
    TYPE5VNIHIGH | auto_evpn_fidevid2vni(fid, evi, 0)
}

```

Figure 13: auto-evpn_fidevi2type5vni

```

/// type-2 VNI for a specific VLAN
pub fn auto_evpn_fidevivid2vni(fid: FabricIDType, evi: MACVRFNumberType, vlanid:
VLANIDType) -> VNIDType {
    let rfid = fid as i32;
    let revi = evi as i32;
    let rvlan = vlanid as i32;
    // mask out high bits, VNI is only 24 bits
    TYPE2VNIMASK &
        (
            rfid.rotate_left(16) ^
            revi.rotate_left(12) ^
            rvlan
        )
}

```

Figure 14: auto_evpn_fidevivid2vni

```

/// maximum VLANs per EVI supported by auto evpn when deriving
pub fn auto_evpn_vlan_description_table<'a>(vlans: usize)
    -> Result<'a [(VLANIDType, VLANStret
chableType, VLANNativeType)], ServiceErrorType> {
    // up to 15 vlans can be activated
    const VLANSARRAY: [(i16, bool, bool); MAX_VLANS_PER_EVI] = [
        (NATIVE_VLAN, true, true, ),
        (FIRST_VLAN_SHIFT, true, false, ),
        (FIRST_VLAN_SHIFT + 1, true, false, ),
        (FIRST_VLAN_SHIFT + 2, true, false, ),
        (FIRST_VLAN_SHIFT + 3, false, false, ),
        (FIRST_VLAN_SHIFT + 4, false, false, ),
        (FIRST_VLAN_SHIFT + 5, false, false, ),
        (FIRST_VLAN_SHIFT + 6, false, false, ),
        (FIRST_VLAN_SHIFT + 7, false, false, ),
        (FIRST_VLAN_SHIFT + 8, false, false, ),
        (FIRST_VLAN_SHIFT + 9, false, false, ),
        (FIRST_VLAN_SHIFT + 10, false, false, ),
        (FIRST_VLAN_SHIFT + 11, false, false, ),
        (FIRST_VLAN_SHIFT + 12, false, false, ),
        (FIRST_VLAN_SHIFT + 13, false, false, ),
    ];

    if vlans > VLANSARRAY.len() {
        return Err(ServiceErrorType::INVALIDPARAMETERVALUE)
    }

    Ok(&VLANSARRAY[..vlans])
}

```

Figure 15: auto_evpn_vlan_description_table

```

/// delivers the vlan description that can be used to generate vlans for a
/// specific fabric ID and a MACVRF number
pub fn auto_evpn_fidevivlansvlans2desc(fid: FabricIDType, macvrf: MACVRFNumberType,
e,
                                vlans: usize) -> Vec<VLANDescription> {

    assert!(NO_MACVRF != macvrf);

    // abstract description of derived VLANs
    let vlan_table = auto_evpn_vlan_description_table(vlans)
        .expect("vlan table in AUTO EVPN incorrect");

    let vlanshift = vlan_table
        .iter()
        .map(|(vl, _, _)| *vl as usize)
        .max()
        .expect("vlan table in AUTO EVPN incorrect")
        .checked_next_power_of_two()
        .expect("vlan table in AUTO EVPN incorrect");

    assert!(vlan_table.len() < FIRST_VLAN_SHIFT as _);

    vlan_table
        .iter()
        .map(move |(vid, stretch, native_)| {
            let stretchedfid = if !stretch {
                fid
            } else {
                NO_FABRIC_ID
            };

            let mut vlan_id = *vid ^ stretchedfid
                .rotate_left(max(16, vlanshift as u32 + 8)) as VLANIDType;
            // leave space for VLANs in the encoding
            vlan_id ^= macvrf.rotate_left(vlanshift as _) as VLANIDType;

            vlan_id %= MAX_VLAN;
            vlan_id = max(1, vlan_id);

            VLANDescription {
                vlan_id: vlan_id as _,
                name: format!("V{}", vlan_id),
                stretchable: *stretch,
                native: *native_,
            }
        })
        .collect()
}

```

Figure 16: auto_evpn_fidevivlansvlans2desc

```

/// IRB interface number.
/// fid/evi combination shifted up to not interfere with the VLAN-ID
/// and then add the VLAN-ID
pub fn auto_evpn_fidevivid2irb(fid: FabricIDType, evi: MACVRFNumberType, vid: VLANIDType) -> UnitType {

    assert!(NO_MACVRF != evi);

    let mut v = (fid as UnitType ^ evi.rotate_left(4) as UnitType) << (16 - FIRST_VLAN.leading_zeros());

    v = 1 + v.wrapping_add(vid) % MAX_VLAN;
    v % (UnitType::MAX - 1)
}

```

Figure 17: auto_evpn_fidevivid2irb

```

/// route distinguisher derivation
pub fn auto_evpn_sidfid2rd(sid: UnsignedSystemID, fid: FabricIDType, extra: u32)
-> RTDistinguisherType {
    // generate type 0 route distinguisher, first 2 bytes 0 and then 6 bytes
    assert!(fid != NO_FABRIC_ID);
    // shift the 2 bytes we loose
    let convsid = sid as RTDistinguisherType;
    let hs = ((sid & 0xffff_0000_0000_0000) >> 32) as RTDistinguisherType;
    let mut ls: RTDistinguisherType = convsid & 0x0000_ffff_ffff_ffff;
    ls ^= hs;
    ls ^= (fid as RTDistinguisherType).rotate_left(16);
    ls ^= extra as RTDistinguisherType;
    ls
}

```

Figure 18: auto_evpn_sidfid2rd


```

/// v4 subnet derivation
pub fn auto_evpn_v4prefixfidevividsid2v4subnet(v4pref: &str, fid: FabricIDType,
                                              evi: MACVRFNumberType, vid: VLANID
Type,
                                              sid: UnsignedSystemID) -> Result<I
Pv4PrefixType, ServiceErrorType> {

    assert!(NO_MACVRF != evi);

    // fid can be 0 for stretched v4subnets
    let mut sub = evi.to_ne_bytes().iter()
        .fold((RANDOMSEEDS[0] & 0xff) as u8, |r, e| r.rotate_left(1) ^ e.rotate_r
ight(1));
    sub ^= fid.to_ne_bytes().iter()
        .fold((RANDOMSEEDS[1] & 0xff) as u8, |r, e| r.rotate_left(2) ^ e.rotate_r
ight(1));
    sub ^= vid.to_ne_bytes().iter()
        .fold((RANDOMSEEDS[2] & 0xff) as u8, |r, e| r.rotate_left(3) ^ e.rotate_r
ight(1));

    let subnet = sub % 254; // make sure we don't show multicast subnet

    let _host = sid.to_ne_bytes().iter()
        .fold(0u16, |r, e| r.rotate_left(3) ^ e.rotate_right(3) as u16);

    let a = format!("{}",
                    v4pref,
                    subnet,
                    0,
                    1,
    );

    Ok(
        IPv4PrefixType {
            address: Ipv4Addr::from_str(&a)
                .map_err(|_| {
                    ServiceErrorType::INTERNALRIFTEERROR
                })?
                .octets()
                .iter()
                .fold(0u32, |v, nv| v << 8 | (*nv as u32)) as Ipv4Address
            ,
            prefixlen: 16,
        }
    )
}

```

Figure 19: auto_evpn_v4prefixfidevividsid2v4subnet

```

/// generic v6 bytes derivation used for different purposes
pub fn auto_evpn_v6hash(fid: FabricIDType, evi: MACVRFNumberType, vid: VLANIDType
, sid: UnsignedSystemID)
    -> [u8; 8] {

    let mut sub = evi.to_ne_bytes().iter()
        .fold(RANDOMSEEDS[3], |r, e| r.rotate_left(6) ^ e.rotate_right(4) as u64)
;
    sub ^= fid.to_ne_bytes().iter()
        .fold(RANDOMSEEDS[0], |r, e| r.rotate_left(6) ^ e.rotate_right(4) as u64)
;
    sub ^= vid as u64;
    sub ^= sid;

    sub.to_ne_bytes()
}

```

Figure 20: auto_evpn_v6hash

```

pub fn auto_evpn_fidevividsid2v6subnet(fid: FabricIDType, evi: MACVRFNumberType,
                                         vid: VLANIDType,
                                         sid: UnsignedSystemID) -> Result<IPv6PrefixType, ServiceErrorType> {

    assert!(NO_MACVRF != evi);

    let sb = auto_evpn_v6hash(fid, evi, vid, sid);

    let a = format!("{:02X}{:02X}{:02X}{:02X}{:02X}{:02X}::1",
        AUTO_EVPN_V6IRBPREF,
        AUTO_EVPN_V6IRBPREFFIRSTBYTE,
        sb[3] ^ sb[0],
        sb[4] ^ sb[1],
        sb[5] ^ sb[2],
        sb[6],
        sb[7],
    );

    Ok(IPv6PrefixType {
        address: Ipv6Addr::from_str(
            &a)
            .map_err(|_| {
                ServiceErrorType::INTERNALRIFTERROR
            })?
            .octets()
            .to_vec(),
        prefixlen: 64,
    })
}

```

Figure 21: auto_evpn_fidevividsid2v6subnet

```

/// MAC address derivation for IRB
pub fn auto_evpn_fidevividssid2mac(fid: FabricIDType, evi: MACVRFNumberType,
                                   vid: VLANIDType, sid: UnsignedSystemID) -> Vec<
u8> {

    let sb = auto_evpn_v6hash(fid, evi, vid, sid);

    vec![0x02,
        sb[3] ^ sb[0],
        sb[4] ^ sb[1],
        sb[5] ^ sb[2],
        sb[6],
        sb[7],
    ]
}

```

Figure 22: auto_evpn_fidevividssid2mac

```

/// v4 loopback address derivation for every node in auto-evpn
pub fn auto_evpn_fidsid2v4loopback(fid: FabricIDType, sid: UnsignedSystemID) -> I
Pv4Address {
    let mut derived = sid.to_ne_bytes().iter()
        .fold(0 as IPv4Address, |p, e| (p << 4) ^ (*e as IPv4Address));
    derived ^= fid as IPv4Address;
    // use the byte we loose for entropy
    derived ^= derived >> 24;
    // and sanitize for loopback range
    derived &= 0x00ff_ffff;

    let m = ((127 as IPv4Address) << 24) | derived;
    m as _
}

```

Figure 23: auto_evpn_fidsid2v4loopback

```

/// V6 loopback derivation for every node in auto-evpn
pub fn auto_evpn_fidsidv6loopback(fid: FabricIDType,
                                   sid: UnsignedSystemID) -> Result<Ipv6Addr, Serv
iceErrorType> {
    auto_evpn_v6prefixfidsid2loopback(AUTO_EVPN_V6PREF, fid, sid)
}

```

Figure 24: auto_evpn_fidsidv6loopback

```

#[allow(non_snake_case)]
pub fn auto_evpn_fid2private_AS(fid: FabricIDType) -> u32 {
    assert!(fid != NO_FABRIC_ID);
    // range 4200000000-4294967294
    const DIFF: u32 = 4_294_967_294 - 4_200_000_000;
    64496 + ((fid as u32) << 3) % DIFF
}

```

Figure 25: auto_evpn_fid2private_AS

```
pub fn auto_evpn_fid2clusterid(fid: FabricIDType) -> u32 {  
    auto_evpn_fid2private_AS(fid)  
}
```

Figure 26: auto_evpn_fid2clusterid

Authors' Addresses

Jordan Head (editor)
Juniper Networks
1137 Innovation Way
Sunnyvale, CA
United States of America

Email: jhead@juniper.net

Tony Przygienda
Juniper Networks
1137 Innovation Way
Sunnyvale, CA
United States of America

Email: prz@juniper.net

Wen Lin
Juniper Networks
10 Technology Park Drive
Westford, MA
United States of America

Email: wlin@juniper.net

RIFT Working Group
Internet-Draft
Intended status: Standards Track
Expires: 12 January 2022

A. Przygienda, Ed.
Juniper
A. Sharma
Comcast
P. Thubert
Cisco
Bruno. Rijsman
Individual
Dmitry. Afanasiev
Yandex
11 July 2021

RIFT: Routing in Fat Trees
draft-ietf-rift-rift-13

Abstract

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

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 12 January 2022.

Copyright Notice

Copyright (c) 2021 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 | 4 |
| 1.1. Requirements Language | 7 |
| 2. A Reader's Digest | 7 |
| 3. Reference Frame | 9 |
| 3.1. Terminology | 9 |
| 3.2. Topology | 15 |
| 4. RIFT: Routing in Fat Trees | 16 |
| 4.1. Overview | 16 |
| 4.1.1. Properties | 17 |
| 4.1.2. Generalized Topology View | 17 |
| 4.1.3. Fallen Leaf Problem | 29 |
| 4.1.4. Discovering Fallen Leaves | 31 |
| 4.1.5. Addressing the Fallen Leaves Problem | 32 |
| 4.2. Specification | 33 |
| 4.2.1. Transport | 34 |
| 4.2.2. Link (Neighbor) Discovery (LIE Exchange) | 35 |
| 4.2.3. Topology Exchange (TIE Exchange) | 49 |
| 4.2.4. Reachability Computation | 74 |
| 4.2.5. Automatic Disaggregation on Link & Node Failures | 76 |
| 4.2.6. Attaching Prefixes | 82 |
| 4.2.7. Optional Zero Touch Provisioning (ZTP) | 90 |
| 4.3. Further Mechanisms | 102 |
| 4.3.1. Route Preferences | 102 |
| 4.3.2. Overload Bit | 103 |
| 4.3.3. Optimized Route Computation on Leaves | 103 |
| 4.3.4. Mobility | 104 |
| 4.3.5. Key/Value Store | 107 |
| 4.3.6. Interactions with BFD | 108 |
| 4.3.7. Fabric Bandwidth Balancing | 109 |
| 4.3.8. Label Binding | 111 |
| 4.3.9. Leaf to Leaf Procedures | 111 |
| 4.3.10. Address Family and Multi Topology Considerations | 112 |
| 4.3.11. One-Hop Healing of Levels with East-West Links | 112 |
| 4.4. Security | 112 |
| 4.4.1. Security Model | 112 |
| 4.4.2. Security Mechanisms | 114 |
| 4.4.3. Security Envelope | 115 |
| 4.4.4. Weak Nonces | 118 |
| 4.4.5. Lifetime | 120 |
| 4.5. Security Association Changes | 120 |

| | | |
|---------|---|-----|
| 5. | Examples | 120 |
| 5.1. | Normal Operation | 120 |
| 5.2. | Leaf Link Failure | 122 |
| 5.3. | Partitioned Fabric | 123 |
| 5.4. | Northbound Partitioned Router and Optional East-West Links | 125 |
| 6. | Further Details on Implementation | 126 |
| 6.1. | Considerations for Leaf-Only Implementation | 126 |
| 6.2. | Considerations for Spine Implementation | 127 |
| 7. | Security Considerations | 127 |
| 7.1. | General | 127 |
| 7.2. | Malformed Packets | 128 |
| 7.3. | ZTP | 128 |
| 7.4. | Lifetime | 128 |
| 7.5. | Packet Number | 128 |
| 7.6. | Outer Fingerprint Attacks | 129 |
| 7.7. | TIE Origin Fingerprint DoS Attacks | 129 |
| 7.8. | Host Implementations | 129 |
| 8. | IANA Considerations | 130 |
| 8.1. | Requested Multicast and Port Numbers | 130 |
| 8.2. | Requested Registries with Suggested Values | 130 |
| 8.2.1. | Registry RIFT_v5/common/AddressFamilyType" | 131 |
| 8.2.2. | Registry RIFT_v5/common/HierarchyIndications" | 131 |
| 8.2.3. | Registry RIFT_v5/common/IEEE802_1ASTimestampType" | 131 |
| 8.2.4. | Registry RIFT_v5/common/IPAddressType" | 132 |
| 8.2.5. | Registry RIFT_v5/common/IPPrefixType" | 132 |
| 8.2.6. | Registry RIFT_v5/common/IPv4PrefixType" | 133 |
| 8.2.7. | Registry RIFT_v5/common/IPv6PrefixType" | 133 |
| 8.2.8. | Registry RIFT_v5/common/PrefixSequenceType" | 133 |
| 8.2.9. | Registry RIFT_v5/common/RouteType" | 134 |
| 8.2.10. | Registry RIFT_v5/common/TIETypeType" | 135 |
| 8.2.11. | Registry RIFT_v5/common/TieDirectionType" | 135 |
| 8.2.12. | Registry RIFT_v5/encoding/Community" | 136 |
| 8.2.13. | Registry RIFT_v5/encoding/KeyValueTIEElement" | 136 |
| 8.2.14. | Registry RIFT_v5/encoding/LIEPacket" | 137 |
| 8.2.15. | Registry RIFT_v5/encoding/LinkCapabilities" | 139 |
| 8.2.16. | Registry RIFT_v5/encoding/LinkIDPair" | 139 |
| 8.2.17. | Registry RIFT_v5/encoding/Neighbor" | 140 |
| 8.2.18. | Registry RIFT_v5/encoding/NodeCapabilities" | 141 |
| 8.2.19. | Registry RIFT_v5/encoding/NodeFlags" | 141 |
| 8.2.20. | Registry RIFT_v5/encoding/NodeNeighborsTIEElement" | 142 |
| 8.2.21. | Registry RIFT_v5/encoding/NodeTIEElement" | 142 |
| 8.2.22. | Registry RIFT_v5/encoding/PacketContent" | 143 |
| 8.2.23. | Registry RIFT_v5/encoding/PacketHeader" | 144 |
| 8.2.24. | Registry RIFT_v5/encoding/PrefixAttributes" | 144 |
| 8.2.25. | Registry RIFT_v5/encoding/PrefixTIEElement" | 145 |
| 8.2.26. | Registry RIFT_v5/encoding/ProtocolPacket" | 145 |
| 8.2.27. | Registry RIFT_v5/encoding/TIDEPacket" | 146 |

| | |
|--|-----|
| 8.2.28. Registry RIFT_v5/encoding/TIEElement" | 146 |
| 8.2.29. Registry RIFT_v5/encoding/TIEHeader" | 147 |
| 8.2.30. Registry RIFT_v5/encoding/TIEHeaderWithLifeTime" | 147 |
| 8.2.31. Registry RIFT_v5/encoding/TIEID" | 148 |
| 8.2.32. Registry RIFT_v5/encoding/TIEPacket" | 148 |
| 8.2.33. Registry RIFT_v5/encoding/TIREPacket" | 149 |
| 9. Acknowledgments | 149 |
| 10. Contributors | 150 |
| 11. References | 150 |
| 11.1. Normative References | 150 |
| 11.2. Informative References | 152 |
| Appendix A. Sequence Number Binary Arithmetic | 154 |
| Appendix B. Information Elements Schema | 155 |
| B.1. Backwards-Compatible Extension of Schema | 156 |
| B.2. common.thrift | 157 |
| B.3. encoding.thrift | 162 |
| Appendix C. Constants | 169 |
| C.1. Configurable Protocol Constants | 169 |
| Authors' Addresses | 171 |

1. Introduction

Clos [CLOS] topologies (called commonly a fat tree/network in modern IP fabric considerations [VAHDATA08] as homonym to the original definition of the term [FATTREE]) 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. Many builders of such IP fabrics desire a protocol that auto-configures itself and deals with failures and mis-configurations with a minimum of human intervention. 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 such IP fabric requirements, RIFT addresses those challenges not through an incremental modification of either a link-state (distributed computation) or distance-vector (diffused computation) techniques but rather a mixture of both, colloquially best described as "link-state towards the spines" 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.

The protocol does further provide

- * optional fully automated construction of fat-tree topologies based on detection of links without any configuration (Section 4.2.7) while allowing for traditional configuration and arbitrary mix of both types of nodes as well,
- * minimum amount of routing state held at each level,
- * automatic pruning and load balancing of topology flooding exchanges over a sufficient subset of links which resolves the traditional problem of link-state protocol struggling with densely meshed graphs due to high volume of flooding traffic (Section 4.2.3.9),
- * automatic aggregation (Section 4.2.3.8) and consequently automatic disaggregation (Section 4.2.5) of prefixes on link and node failures to prevent black-holing and suboptimal routing,
- * loop-free non-ECMP forwarding due to its inherent valley-free nature,
- * fast mobility (Section 4.3.4),
- * re-balancing of traffic towards the spines based on bandwidth available (Section 4.3.7.1) and finally
- * mechanisms to synchronize a limited key-value data-store (Section 4.3.5.1) that can be used after protocol convergence to e.g. bootstrap higher levels of functionality on nodes.

Figure 1 presents as first example of operation a simplified, conceptual view of the resulting information and routes on a RIFT fabric. The top of the fabric is holding in its link-state database the information about the nodes below it and the routes to them whereas the notation A/32 is used to indicate a loopback route to node A and 0/0 is the usual notation for a default route. First row of information represents the nodes for which full topology information is available. The second row of the database table indicates that partial information of other nodes in the same level is available as well. Such information will be necessary to perform certain algorithms necessary for correct protocol operation. When "bottom" of the fabric is considered, or in other words the leaves, the topology is basically empty and, under normal conditions, the leaves hold a load balanced default route to the next level.

The balance of this document fills in the protocol specification details.

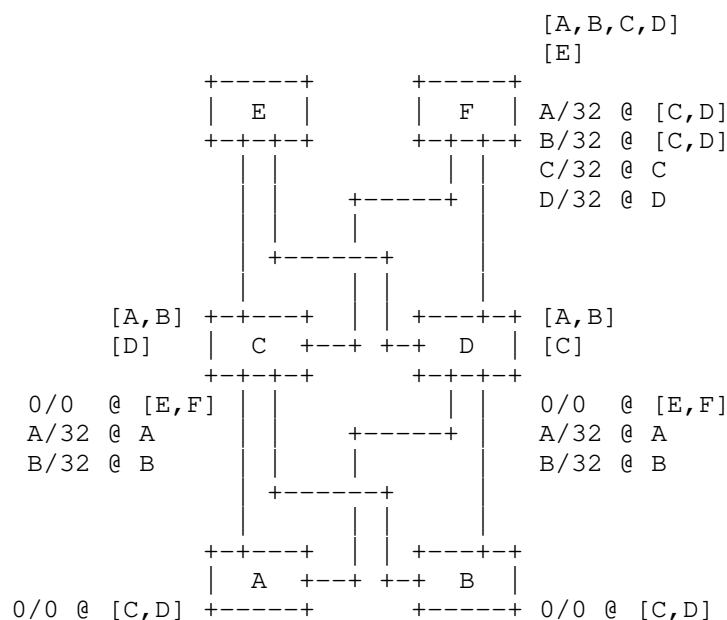


Figure 1: RIFT Information Distribution

1.1. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 RFC 2119 [RFC2119] RFC 8174 [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. A Reader's Digest

This section should serve as an initial guided tour through the document in order to convey the necessary information for any reader, depending on their level of interest. The glossary section (Section 3.1) should be used as a supporting reference as the document is read.

The indications to direction (i.e. "top", "bottom", etc.) referenced in the Section 1 are of paramount importance. RIFT requires a topology with a sense top and bottom in order to properly achieve a sorted topology. Clos, Fat-Tree, and other similarly structured networks are conducive to such requirements. RIFT does allow for further relaxation of these constraints, they will be mentioned later in this section.

Operators and implementors alike must understand if multi-plane IP fabrics are of interest or not. Section 3.2 illustrates an example of both single-plane in Figure 2 and multi-plane fabric in Figure 3. Multi-plane fabrics require understanding of additional RIFT concepts (e.g. negative disaggregation in Section 4.2.5.2) that are otherwise unnecessary in context of strictly single-plane fabrics. Overview (Section 4.1) and Section 4.1.2 aim to provide enough context to determine if multi-plane fabrics are of interest to the reader. The Fallen Leaf part (Section 4.1.3), and additionally Section 4.1.4 and Section 4.1.5 describe further considerations that are specific to multi-plane fabrics.

The fundamental protocol concepts are described starting in the specification part (Section 4.2), but some sub-sections are not quite as relevant unless dealing with implementation of the protocol. The protocol transport (Section 4.2.1) is of particular importance for two reasons. First, it introduces RIFT's packet formats in the form of a normative Thrift model given in Appendix B.3. Second, the Thrift model component is a prelude to understanding the RIFT's inherent security features as defined in the security segment (Section 7). The normative schema defining the Thrift model can be found in both Appendix B.2 and Appendix B.3. Furthermore, while a detailed understanding of Thrift and the models are not required unless implementing RIFT, they may provide additional useful information for other readers.

If implementing RIFT to support multi-plane topologies Section 4.2 should be reviewed in its entirety in conjunction with previously mentioned Thrift schemas. Sections not relevant to single-plane implementations will be noted later in the section. Special attention should be paid to the LIE definitions part (Section 4.2.2) as it not only outlines basic neighbor discovery and adjacency formation, but also provides necessary context for RIFT's ZTP (Section 4.2.7) and mis-cabling detection capabilities that allow it to automatically detect and build the underlay topology with a negligible configuration. These specific capabilities are detailed in Section 4.2.7.

For other readers, the following sections provide a more detailed understanding of the fundamental properties and highlight some additional benefits of RIFT such as link state packet formats, highly efficient flooding, synchronization, loop-free path computation and link-state database maintenance - Section 4.2.3, Section 4.2.3.2, Section 4.2.3.3, Section 4.2.3.4, Section 4.2.3.6, Section 4.2.3.7, Section 4.2.3.8, Section 4.2.4, Section 4.2.4.1, Section 4.2.4.2, Section 4.2.4.3, Section 4.2.4.4. RIFT's unique ability to perform weighted unequal-cost load balancing of traffic across all available links is outlined in Section 4.3.7 with an accompanying example.

Section 4.2.5 is the place where the single-plane vs. multi-plane requirement is explained in more detail. For those interested in single-plane fabrics, only Section 4.2.5.1 is required. For the multi-plane interested reader Section 4.2.5.2, Section 4.2.5.2.1, Section 4.2.5.2.2, and Section 4.2.5.2.3 are also mandatory. Section 4.2.6 is especially important for any multi-plane interested reader as it outlines how the RIB and FIB are built via the disaggregation mechanisms, but also illustrates how they prevent defective routing decisions (e.g. black holes) in both single or multi-plane topologies.

Section 5 contains a set of comprehensive examples that continue to highlight just how efficiently RIFT handles failures by containing impact to only the required set of nodes. It should also help cement some of RIFT's core concepts in the reader's mind.

Last, but not least, RIFT has other optional capabilities. One example is the key-value data-store, which enables RIFT to advertise data post-convergence in order to bootstrap higher levels of functionality (e.g. operational telemetry). Those are covered in Section 4.3 and Section 6.

More information related to RIFT can be found in the "RIFT Applicability" [APPLICABILITY] document, which discusses alternate topologies upon which RIFT may be deployed, use cases where it is applicable, and presents operational considerations that complement this document.

3. Reference Frame

3.1. Terminology

This section presents the terminology used in this document.

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 a 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 the Clos fabric input and output stages are analogous, the fabric can be "folded" to build a "superspine" or top which is called 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 and a leaf may have links to nodes multiple levels higher. RIFT counts levels from top-of-fabric (ToF) numerically down. Level 0 always implies a leaf in RIFT but a leaf does not have to be level 0. Level in RIFT can be configured or automatically derive its level via ZTP as explained in Section 4.2.7. As final footnote: Clos terminology uses often the concept of "stage" but due to the folded nature of the Fat Tree it is not used from this point on to prevent misunderstandings.

Superspine, Aggregation/Spine and Edge/Leaf Switches:"

Traditional level names in 5-stages folded Clos for Level 2, 1 and 0 respectively (counting up from the bottom). 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. Such a minimum configuration consists solely of ToFs being configured as such.

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. PoDs are numbered to distinguish them and PoD value 0 (defined later in the encoding schema as 'common.default_pod') is used to denote "undefined" or "any" 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 'common.default_pod' 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. As mentioned before, Level 0 implies a leaf in RIFT but a leaf does not have to be level 0.

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. 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 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, in different situations the directionality of the compass is used. I.e., 'lower', 'south' or 'southbound' mean moving towards the bottom of the Clos or Fat Tree network and 'higher', '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 (E-W) 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 Topology Information Elements (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. A TIE has always a direction and a type. North TIEs (sometimes abbreviated as N-TIEs) are used when dealing with 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", which contains all adjacencies the node discovered and information about the node itself. Node TIE should NOT be confused with a North TIE since "node" defines the type of TIE rather than its direction. Consequently North Node TIEs and South Node TIEs exist.

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 (KV) TIE:

A TIE that is carrying a set of key value pairs [DYNAMO]. It can be used to distribute non topology related information within the protocol.

TIDE:

Topology Information Description Element carrying descriptors of the TIEs stored in the node.

TIRE:

Topology Information Request Element carrying set of TIDE descriptors. It can both confirm received and request missing TIEs.

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" exchanged on all the system's links running RIFT to form ThreeWay adjacencies and carry information used to perform Zero Touch Provisioning (ZTP) of levels.

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. 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 the 'overload' attribute as set. Overload attribute is carried in the 'NodeFlags' object of the encoding schema.

Interface:

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

ThreeWay 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 ThreeWay state, i.e. both routers reflected each other in LIEs including relevant security information. Nevertheless, LIEs before ThreeWay 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 ThreeWay 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 ThreeWay 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 [DIJKSTRA] that establishes a tree of shortest paths from a source to destinations on the graph. SPF acronym is used due to its familiarity as general term for the node reachability calculations RIFT can employ to ultimately calculate routes of which Dijkstra algorithm is a possible 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.

System ID:

Each RIFT node identifies itself by a valid, network wide unique number when trying to build adjacencies or describing its topology. RIFT System IDs can be auto-derived or configured.

Additionally, when the specification refers to elements of packet encoding or constants provided in the Appendix B grave accents are used, e.g. `'invalid_distance'`. Same convention is used when referring to finite state machine states or events outside the context of the machine itself, e.g. `'OneWay'`.

3.2. Topology

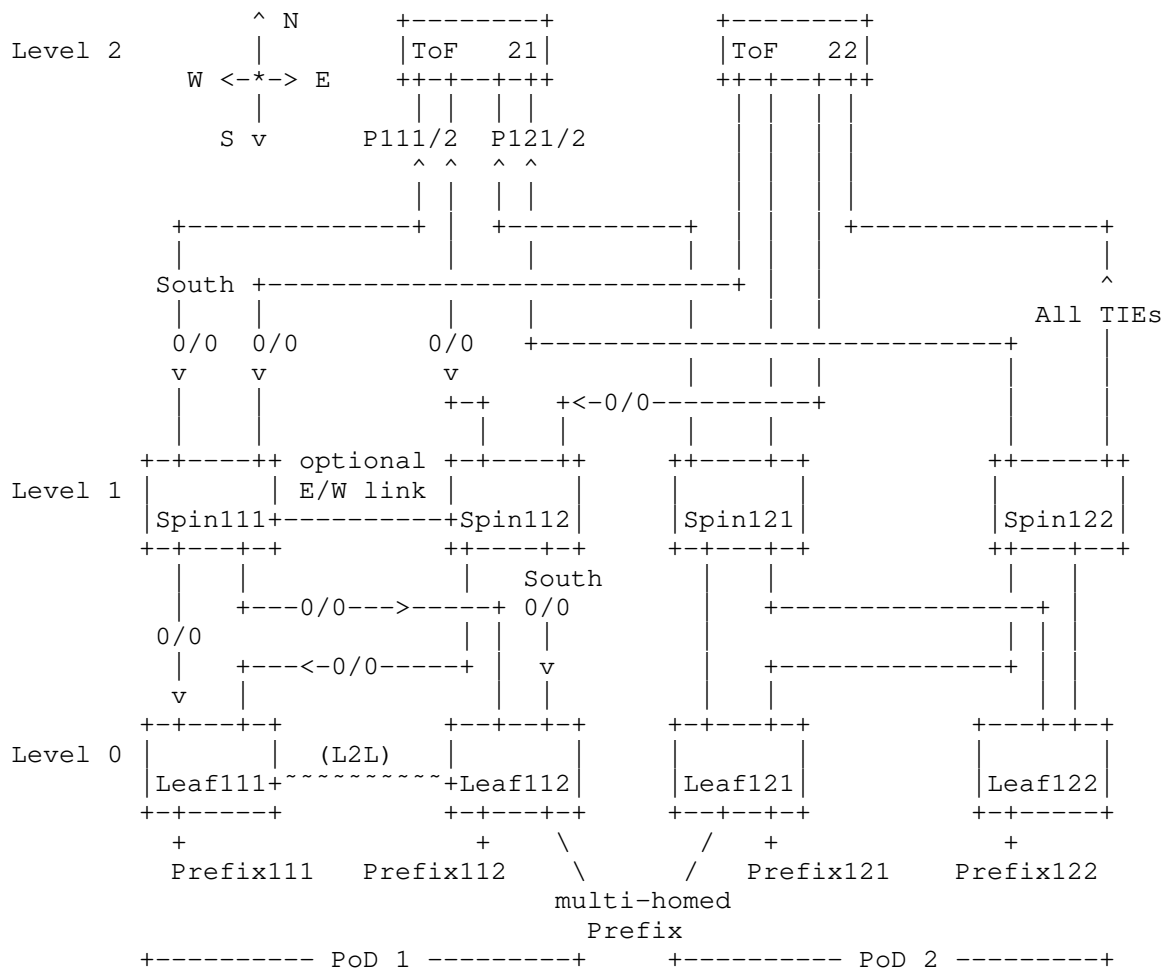


Figure 2: A Three Level Spine-and-Leaf Topology

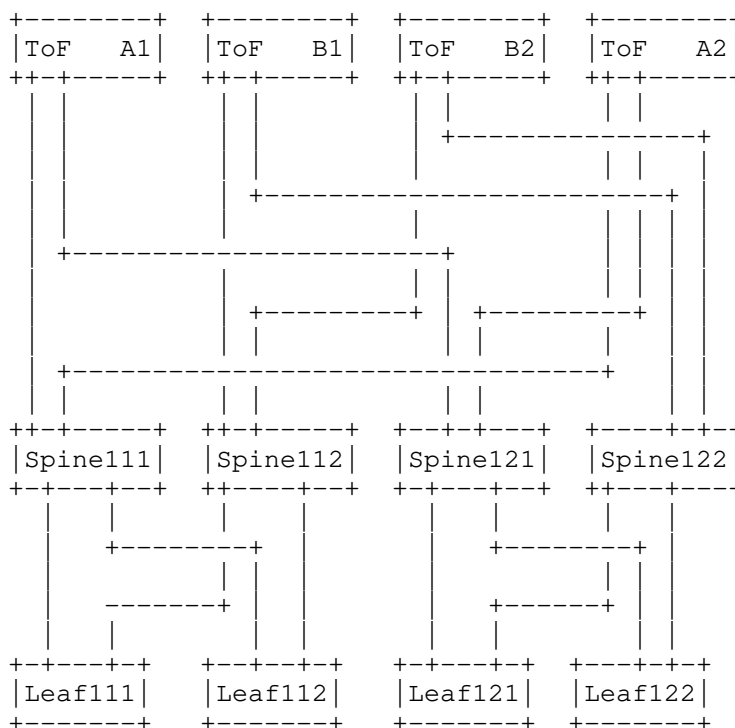


Figure 3: Topology with Multiple Planes

Topology in Figure 2 is referred to 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

Remainder of this documents presents the detailed specification of a protocol optimized for Routing in Fat Trees (RIFT) that in most abstract terms has many properties of a modified link-state protocol when distributing information northbound and a distance vector protocol when distributing information southbound. While this is an unusual combination, it does quite naturally exhibit the desirable properties desired.

4.1. Overview

4.1.1. Properties

The most singular property of RIFT is that it floods 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 the southbound direction, the necessary 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. For the moment describing the East-West direction is left out.

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 the same reachability information from multiple directions. Its computation principles (south forwarding direction is always preferred) leads to valley-free [VFR] 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 South TIEs carry simply node description of the level above and default routes pointing north. This oversimplified view will be refined gradually in the following sections while introducing protocol procedures and state machines at the same time.

4.1.2. Generalized Topology View

This section and resulting Section 4.2.5.2 are dedicated to multi-plane fabrics, in contrast with the single plane designs where all top-of-fabric nodes are topologically equal and initially connected to all the switches at the level below them.

It is quite difficult to visualize multi plane design, which are effectively multi-dimensional switching matrices. To cope with that, this document introduces a methodology allowing to depict the connectivity in two-dimensional pictures. Further, the fact can be leveraged that what is under consideration here are basically 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 and Glossary

This section describes the terminology and acronyms used in the rest of the text. Though the glossary may not be comprehensible on a first read, the following sections will gradually introduce the terms in their proper context.

P:

Denotes the number of PoDs in a topology.

S:

Denotes the number of ToF nodes in a topology.

K:

To simplify the visual aids, notations and further considerations, implicit assumption is made that the switches are symmetrical, i.e. equal number ports point northbound and southbound. With that simplification, K denotes half of the radix of a symmetrical switch, meaning that the switch has K ports pointing north and K ports pointing south. K_LEAF (K of a leaf) thus represents both the number of access ports in a leaf Node and the maximum number of planes in the fabric, whereas K_TOP (K of a ToP) represents the number of leaves in the PoD and the number of ports pointing north in a ToP Node towards a higher spine level, thus the number of ToF nodes in a plane.

ToF Plane:

Set of ToFs that are aware of each other by means of south reflection. Planes are numbered by capital letters, e.g. plane A.

N:

Denotes 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.

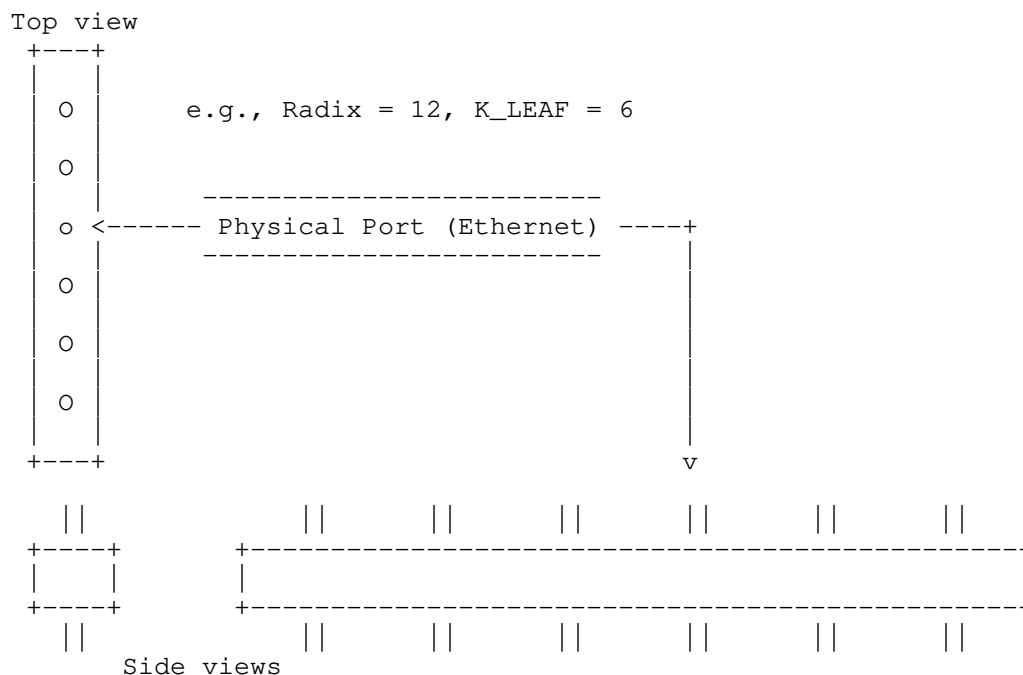
Fallen Leaf:

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

4.1.2.2. Clos as Crossed, Stacked 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. From here on half of them ($K = \text{Radix}/2$) are assumed to connect host devices from the south, and the other half to connect to interleaved PoD Top-Level switches to the north. The K ratio can be chosen differently without loss of generality when port speeds differ or the fabric is oversubscribed but $K = \text{Radix}/2$ allows for more readable representation whereby there are as many ports facing north as south on any intermediate node. A node is hence represented 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, the document chooses 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 ($K \times K$). In the general case, switches at the top of the PoD with K_{TOP} southern ports not necessarily equal to K_{LEAF} could be considered. 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, K_{TOP} Leaf Nodes are necessary as illustrated in Figure 5.

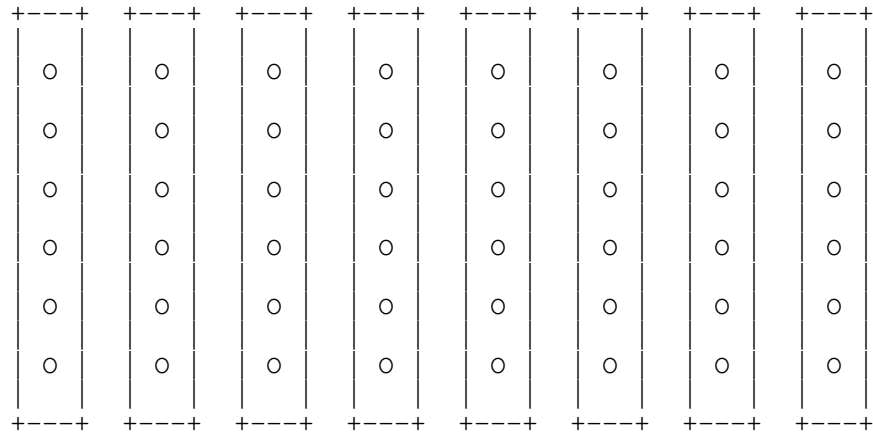


Figure 5: Southern View of a PoD, K_TOP=8

As further visualized in Figure 6 the K_TOP Leaf Nodes are fully interconnected with the K_LEAF 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".

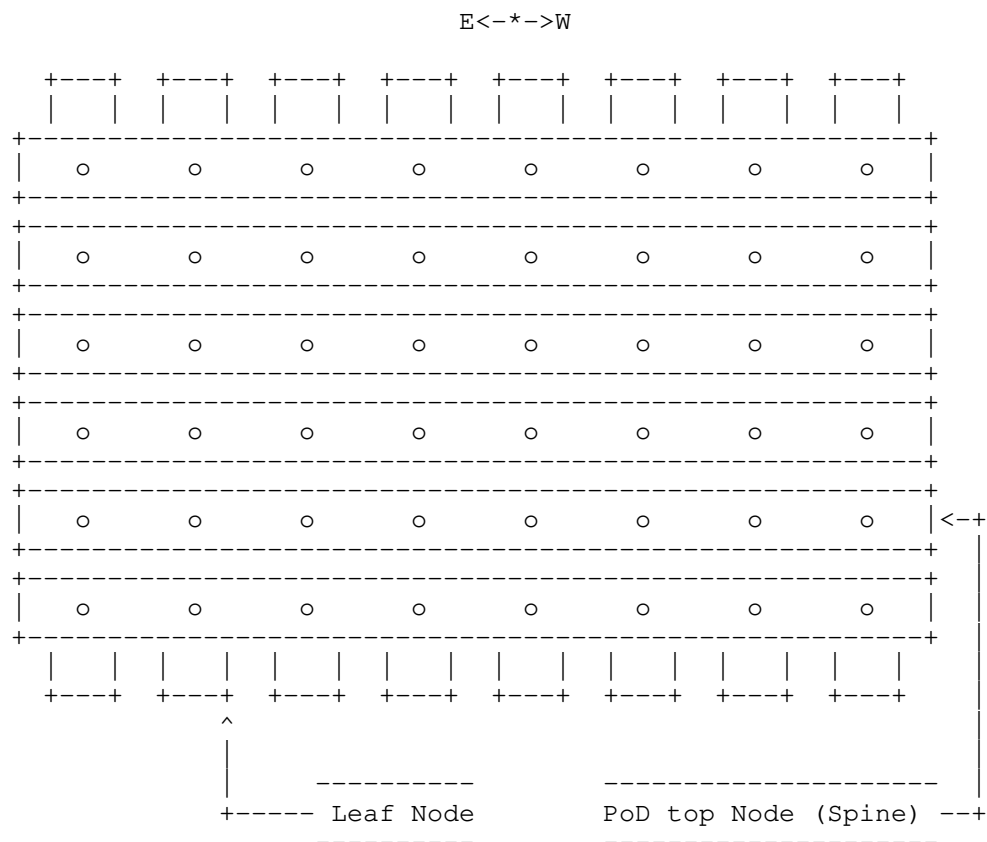


Figure 6: Northern View of a PoD's Spines, K_TOP=8

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

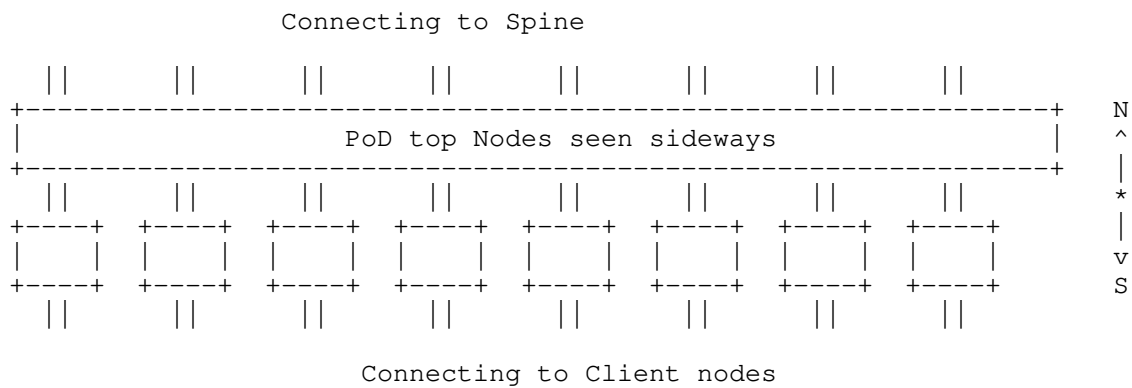
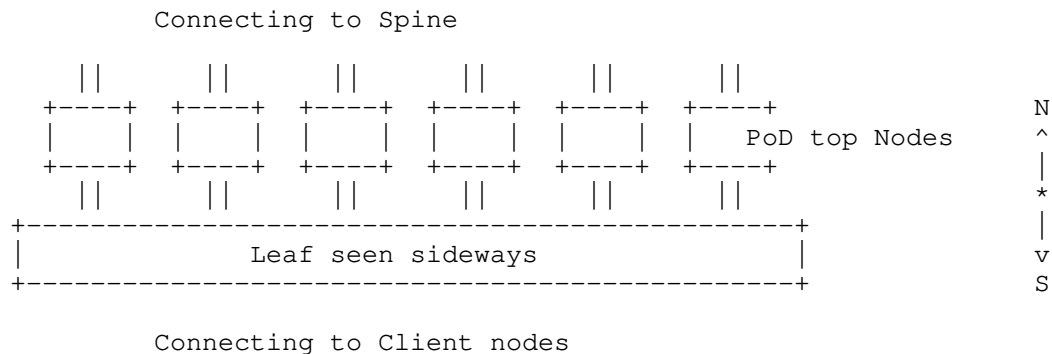


Figure 7: Side View of a PoD, $K_{TOP}=8$, $K_{LEAF}=6$ Figure 8: Other Side View of a PoD, $K_{TOP}=8$, $K_{LEAF}=6$, 90o turn in E-W Plane from the previous figure

As next step, observe further 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. Figure 9 illustrates this, looking at $P=3$ PoDs from above and 2 sides. The large view is the one from above, with the 8 ToF of $3*6$ ports each interconnecting the PoDs, every ToP Node being connected to every ToF node.

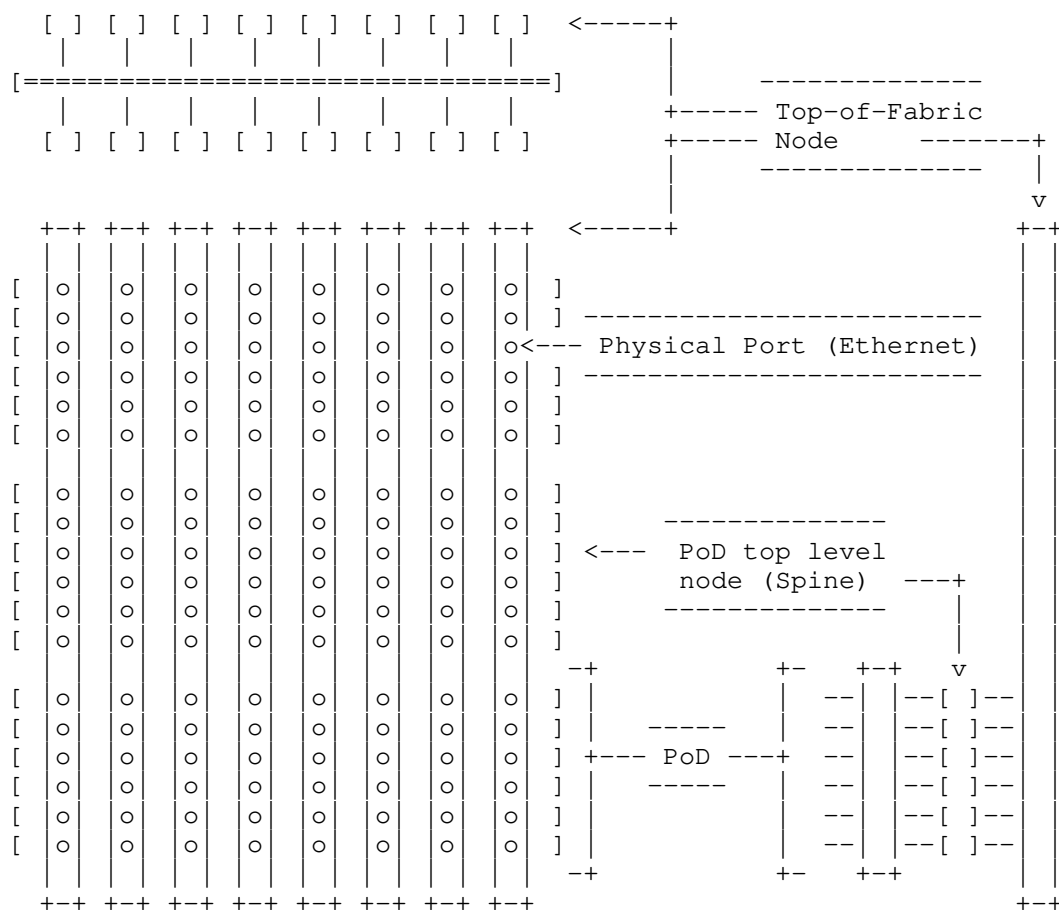


Figure 9: Fabric Spines and TOFs in Single Plane Design, 3 PoDs

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

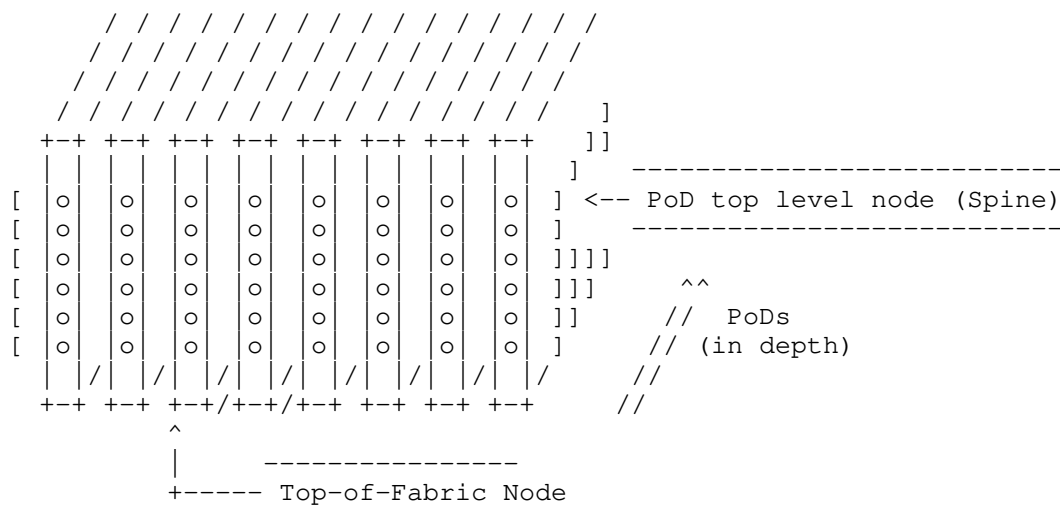


Figure 10: Collapsed Northern View of a Fabric for Any Number of PoDs

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, it will be come clear 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 an 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 that denotes the number of independent paths between 2 leaves within a plane. 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 ToP nodes (spines) in each PoD. The ToP nodes are represented horizontally with $K_{TOP}=8$ ports northwards each.

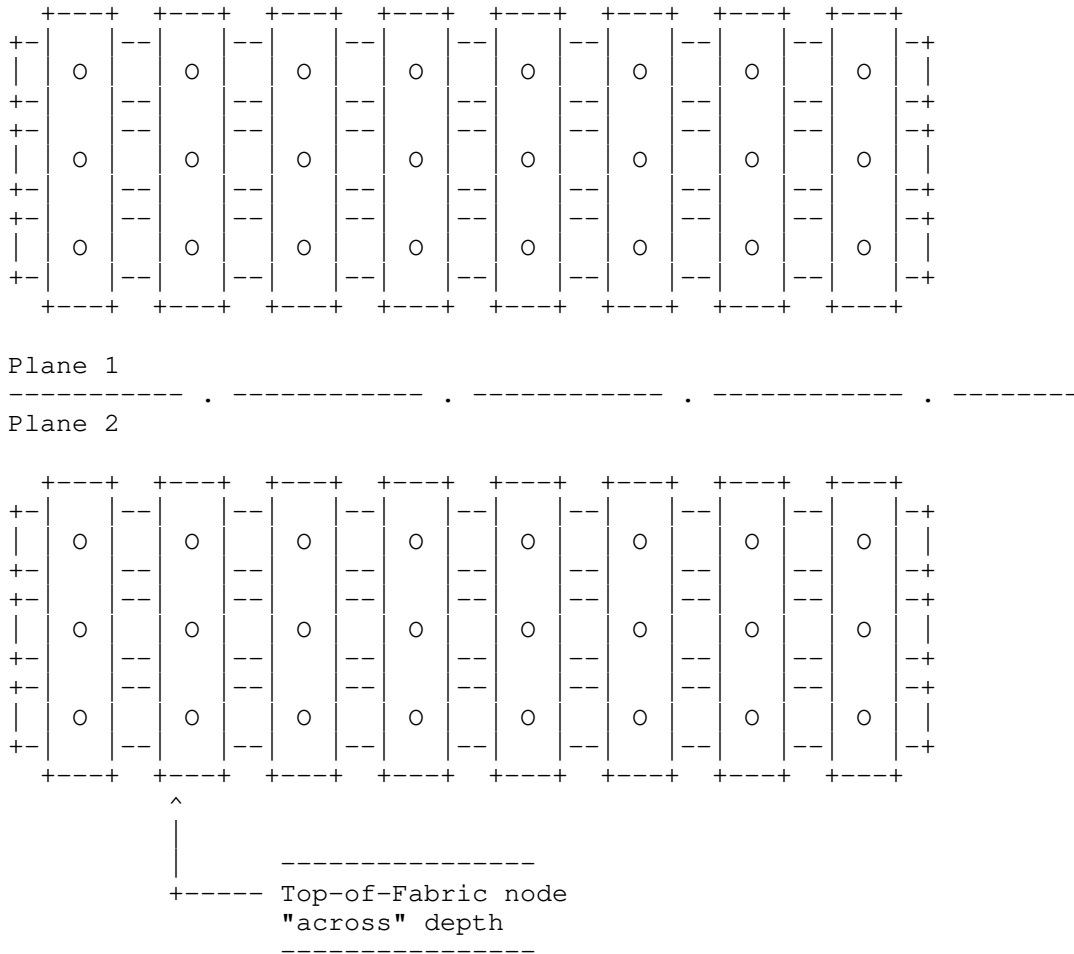


Figure 11: Northern View of a Multi-Plane ToF Level, $K_{LEAF}=6$, $N=2$

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 than or equal to P , the number of PoDs.

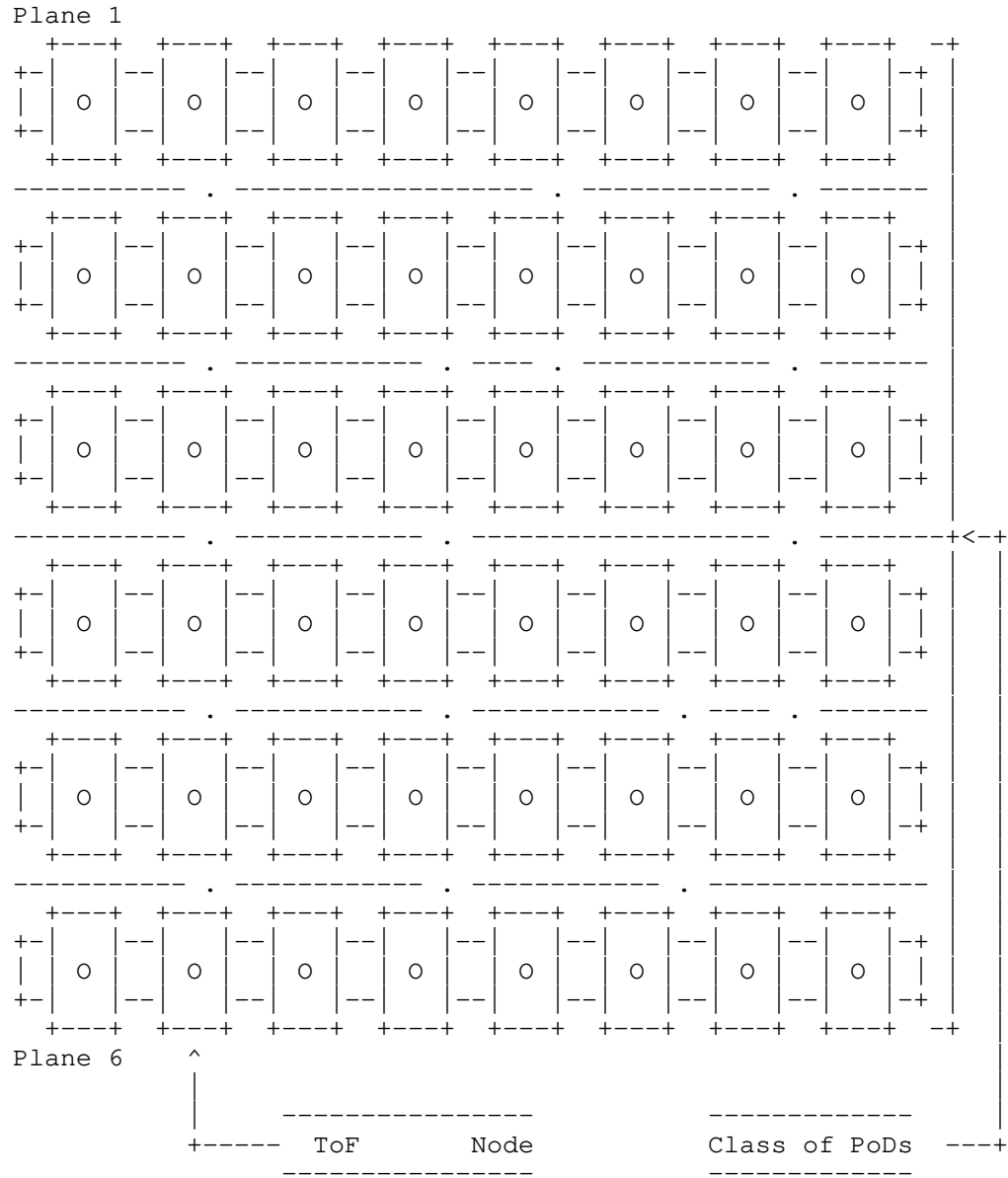


Figure 12: Northern View of a Maximally Partitioned ToF Level, R=1

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.

A "Fallen Leaf" is 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 will cause one or more fallen leaves in the affected plane. $R=2$ guarantees that a single breakage will not cause a fallen leaf. However, not all cases require disaggregation. The following cases do not require particular action:

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 within the plane in a transitive fashion, that is, all the way to the other leaves, in the following cases:

- * 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 the control plane that it cannot reach a certain more specific prefix than default whereas such disaggregation must in the 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.
- * 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 that 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 the abstractions are rolled back into a simple example that shows 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 the 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:

- * 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.
- * A path computation trying to deal with the problem by distributing host routes may only form paths through leaves. The flooding of information about 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

When aggregation is used, RIFT deals with fallen leaves by ensuring that all the ToF nodes 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 enable routing to fallen leaves in multi-plane designs, RIFT requires additional interconnection across planes between the ToF nodes, e.g., using rings as illustrated in Figure 13. 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 requirement to have a full view of the fabric topology at the Top-of-Fabric level, without the need to collate it from multiple points.

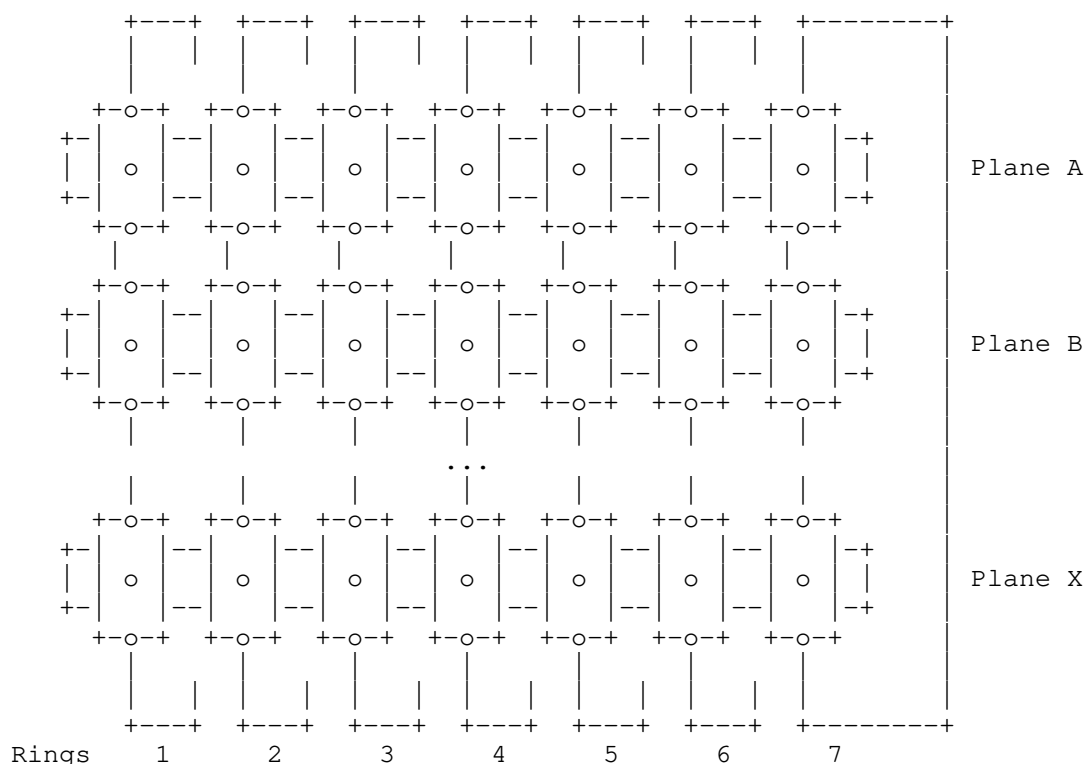


Figure 13: Using rings to bring all planes and at the ToF bind them

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 defines two methods to address this issue, the positive and the negative disaggregation. Both methods flood according types of South TIEs to advertise the impacted prefix(es).

When used for the operation of disaggregation, a positive South TIE contained in `'positive_disaggregation_prefixes'`, as usual, indicates reachability to a prefix of given length and all addresses subsumed by it. In contrast, a negative route advertisement contained in `'negative_disaggregation_prefixes'` 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 TIEs or floods them 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. The specific procedures will be explained in detail in Section 4.2.5.2.3.

When the top of fabric switches are not partitioned into multiple planes, the resulting southbound 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).

The FSMs, as usual, are presented as states the FSM can assume, events that it can be given and according actions performed when transitioning between states on event processing.

Actions are performed before the end state is assumed.

The FSMs can queue events against itself to chain actions or against other FSMs in the specification. Events are always processed in the sequence they have been queued.

Consequently, "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, i.e. the protocol MUST reset all adjacencies, discard all the state and MAY NOT start again.

The data structures and FSMs described in this document are conceptual and do not have to be implemented precisely as described here, as long as the implementations support the described functionality and exhibit the same externally visible behavior.

The machines can use conceptually "timers" for different situations. Those timers are started through actions and their expiration leads to queuing of according events to be processed.

The term 'holdtime' is used often as short-hand for 'holddown timer' and signifies either the length of the holding down period or the timer used to expire after such period. Such timers are used to "hold down" state within an FSM that is cleaned if the machine triggers a 'HoldtimeExpired' event.

4.2.1. Transport

All packet formats are defined in Thrift [thrift] models in Appendix B. LIE packet format is contained in the 'LIEPacket' schema element. TIE packet format is contained in 'TIEPacket', TIDE and TIRE accordingly in 'TIDEPacket', 'TIREPacket' and the whole packet is a union of the above in 'ProtocolPacket' while it contains a 'PacketHeader' as well.

Such a packet being in terms of bits on the wire a serialized 'ProtocolPacket' is carried in an envelope defined in Section 4.4.3 within a UDP frame that provides security and allows validation/modification of several important fields without de-serialization for performance and security reasons. Security model and procedures are further explained in Section 7.

4.2.2. Link (Neighbor) Discovery (LIE Exchange)

RIFT LIE exchange auto-discovers neighbors, negotiates ZTP parameters and discovers miscablings. The formation progresses under normal conditions from OneWay to TwoWay and then ThreeWay state at which point it is ready to exchange TIEs per Section 4.2.3. The adjacency exchanges ZTP information (Section 4.2.7) in any of the states, i.e. it is not necessary to reach ThreeWay for zero-touch provisioning to operate.

RIFT supports any combination of IPv4 and IPv6 addressing on the fabric with the additional capability for forwarding paths that are capable of forwarding IPv4 packets in presence of IPv6 addressing only.

For IPv4 LIE exchange happens over well-known administratively locally scoped and configured or otherwise well-known IPv4 multicast address [RFC2365]. For IPv6 [RFC8200] exchange is performed over link-local multicast scope [RFC4291] address which is configured or otherwise well-known. In both cases a destination UDP port defined in Appendix C.1 is used unless configured otherwise. LIEs SHOULD be sent with an IPv4 Time to Live (TTL) / IPv6 Hop Limit (HL) of either 1 or 255 to prevent RIFT information reaching beyond a single L3 next-hop in the topology. LIEs SHOULD be sent with network control precedence unless an implementation is prevented from doing so.

The 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 LIE FSM 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 TIEs (Section 4.2.3.3). That implies that an implementation MUST be ready to accept TIEs on all addresses it used as source of LIE frames.

A ThreeWay adjacency (as defined in the glossary) over any address family implies support for IPv4 forwarding if the 'ipv4_forwarding_capable' flag in 'LinkCapabilities' is set to true. A node, in case of absence of IPv4 addresses on such links and advertising 'ipv4_forwarding_capable' as true, MUST forward IPv4

packets using gateways discovered on IPv6-only links advertising this capability. It is expected that the whole fabric supports the same type of forwarding of address families on all the links, any other combination is outside the scope of this specification.

'ipv4_forwarding_capable' MUST be set to true when LIEs from a IPv4 address are sent and MAY be set to true in LIEs on IPv6 address if no LIEs are sent from a IPv4 address. If IPv4 and IPv6 LIEs indicate contradicting information protocol behavior is unspecified.

Operation of a fabric where only some of the links are supporting forwarding on an address family or have an address in a family and others do not is outside the scope of this specification.

Any attempt to construct IPv6 forwarding over IPv4 only adjacencies is outside this specification.

Table 1 outlines protocol behavior in case of different address family combinations.

| AF | AF | Behavior |
|---------------|---------------|--|
| IPv4 | IPv4 | LIEs and TIEs are exchanged over IPv4, no IPv6 forwarding. TIEs are received on any of the LIE sending addresses. |
| IPv6 | IPv6 | LIEs and TIEs are exchanged over IPv6 only, no IPv4 forwarding if either of the 'ipv4_forwarding_capable' flags is false. If both 'ipv4_forwarding_capable' flags are true IPv4 is forwarded. TIEs are received on any of the LIE sending addresses. |
| IPv4, IPv6 | IPv6 | LIEs and TIEs are exchanged over IPv6, no IPv4 forwarding if either of the 'ipv4_forwarding_capable' flags is false. If both 'ipv4_forwarding_capable' are true IPv4 is forwarded. TIEs are received on any of the IPv6 LIE sending addresses. |
| IPv4, IPv6 | IPv4, IPv6 | LIEs and TIEs are exchanged over IPv6 and IPv4, unspecified behavior if either of the 'ipv4_forwarding_capable' flags is false or IPv4 and IPv6 advertise different flags as described previously. IPv4 and IPv6 are forwarded. TIEs are received on any of the IPv4 and IPv6 LIE sending addresses. |

Table 1: Neighbor AF Combination Behavior

The protocol does *not* support selective disabling of address families after adjacency formation, disabling IPv4 forwarding capability or any local address changes in ThreeWay state, i.e. if a link has entered ThreeWay 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 IPv4 forwarding, it has to tear down and rebuild the adjacency. It also has to remove any state it stored about the remote side of 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 and advertises it in the 'level' of the 'PacketHeader' schema element. It MAY be also provisioned with its PoD. If level is not provisioned it is not present in the optional 'PacketHeader' schema element and established by ZTP procedures if feasible. If PoD is not provisioned it is as

governed by the 'LIEPacket' schema element assuming the 'common.default_pod' value. This means that switches except top of fabric do not need to be configured at all. Necessary information to configure all values is exchanged in the 'LIEPacket' and 'PacketHeader' or derived by the node automatically.

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

A node MUST form a ThreeWay adjacency (or in other words consider the neighbor "valid" and hence reflecting it) if and only if the following first order logic conditions are satisfied on a LIE packet as specified by the 'LIEPacket' schema element and received on a link

1. the neighboring node is running the same major schema version as indicated in the 'major_version' element in 'PacketHeader' *and*
2. the neighboring node uses a valid System ID (i.e. value different from 'IllegalSystemID') in 'sender' element in 'PacketHeader' *and*
3. the neighboring node uses a different System ID than the node itself
4. the advertised MTUs in 'LiePacket' element match on both sides *and*
5. both nodes advertise defined level values in 'level' element in 'PacketHeader' *and*
6. [
 - i) the node is at 'leaf_level' value and has no ThreeWay adjacencies already to nodes at Highest Adjacency ThreeWay (HAT as defined later in Section 4.2.7.1) with level different than the adjacent node *or*
 - ii) the node is not at 'leaf_level' value and the neighboring node is at 'leaf_level' value *or*
 - iii) both nodes are at 'leaf_level' values *and* both indicate support for Section 4.3.9 *or*
 - iv) neither node is at 'leaf_level' value and the neighboring node is at most one level difference away].

LIEs arriving with IPv4 Time to Live (TTL) / IPv6 Hop Limit (HL) different than 1 or 255 SHOULD be ignored.

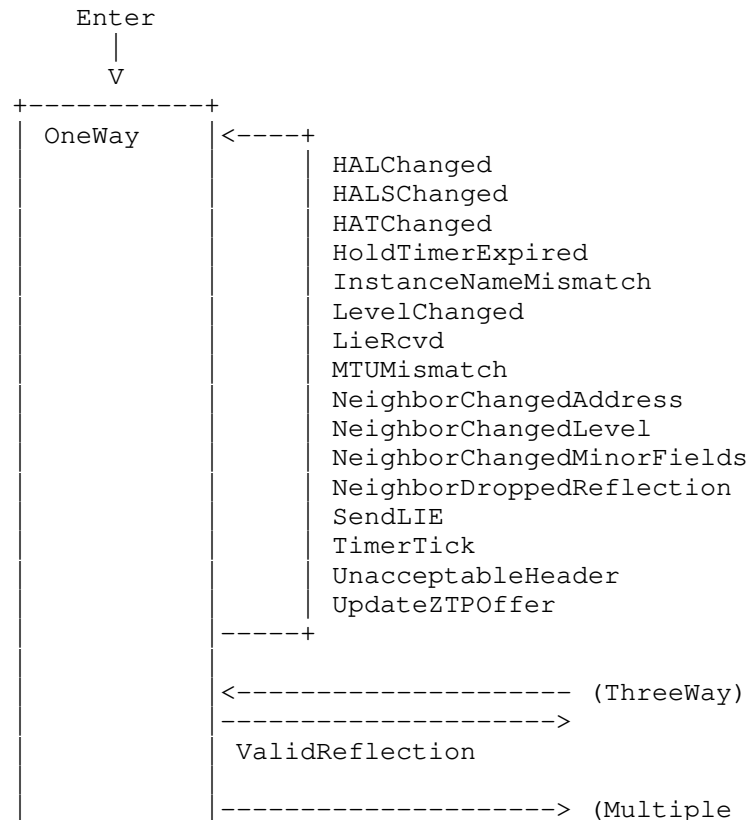
4.2.2.1. LIE Finite State Machine

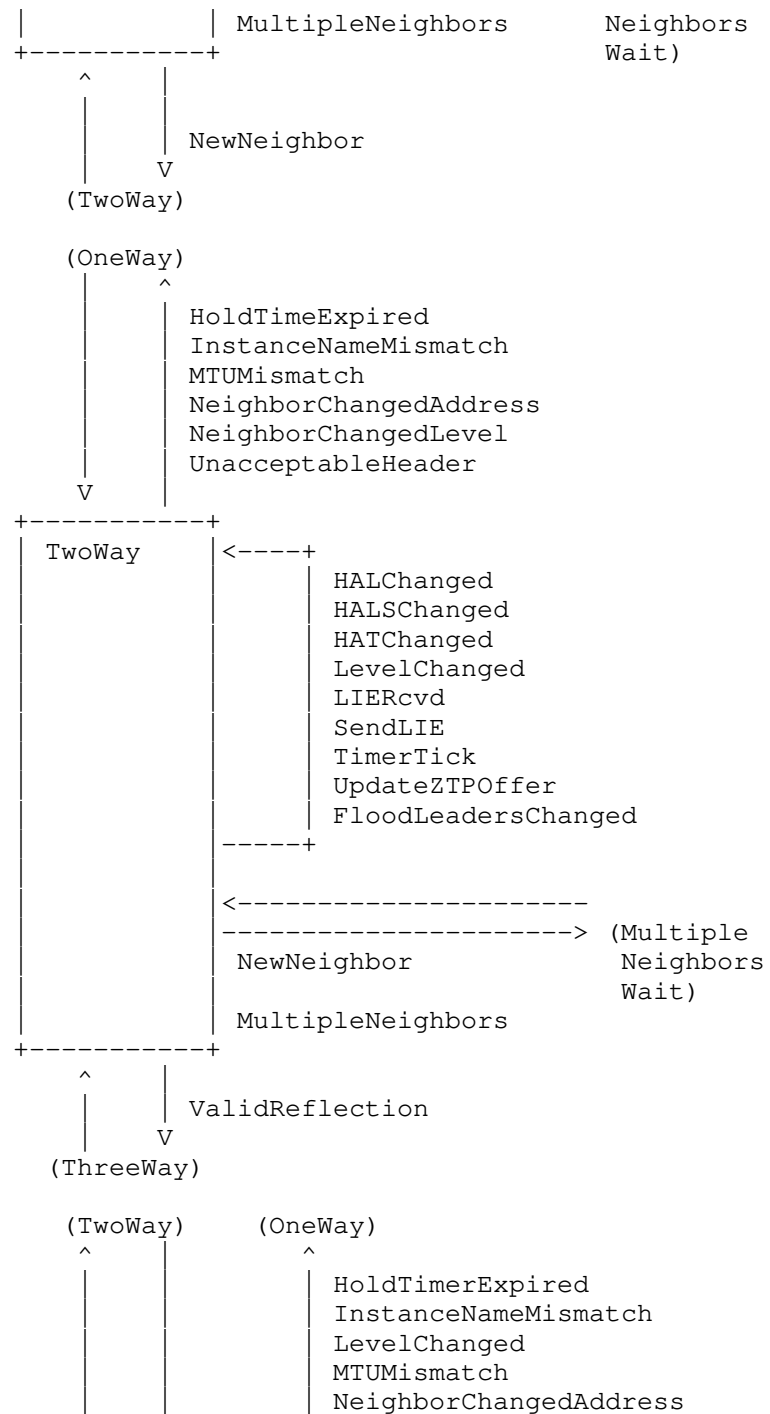
This section specifies the precise, normative LIE FSM. For easier reference the according figure is given as well in Figure 14. Additionally, some sets of actions repeat often and are hence summarized into well-known procedures.

Events generated are fairly fine grained, especially when indicating problems in adjacency forming conditions. The intention of such differentiation is to simplify tracking of problems in deployment.

Initial state is 'OneWay'.

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





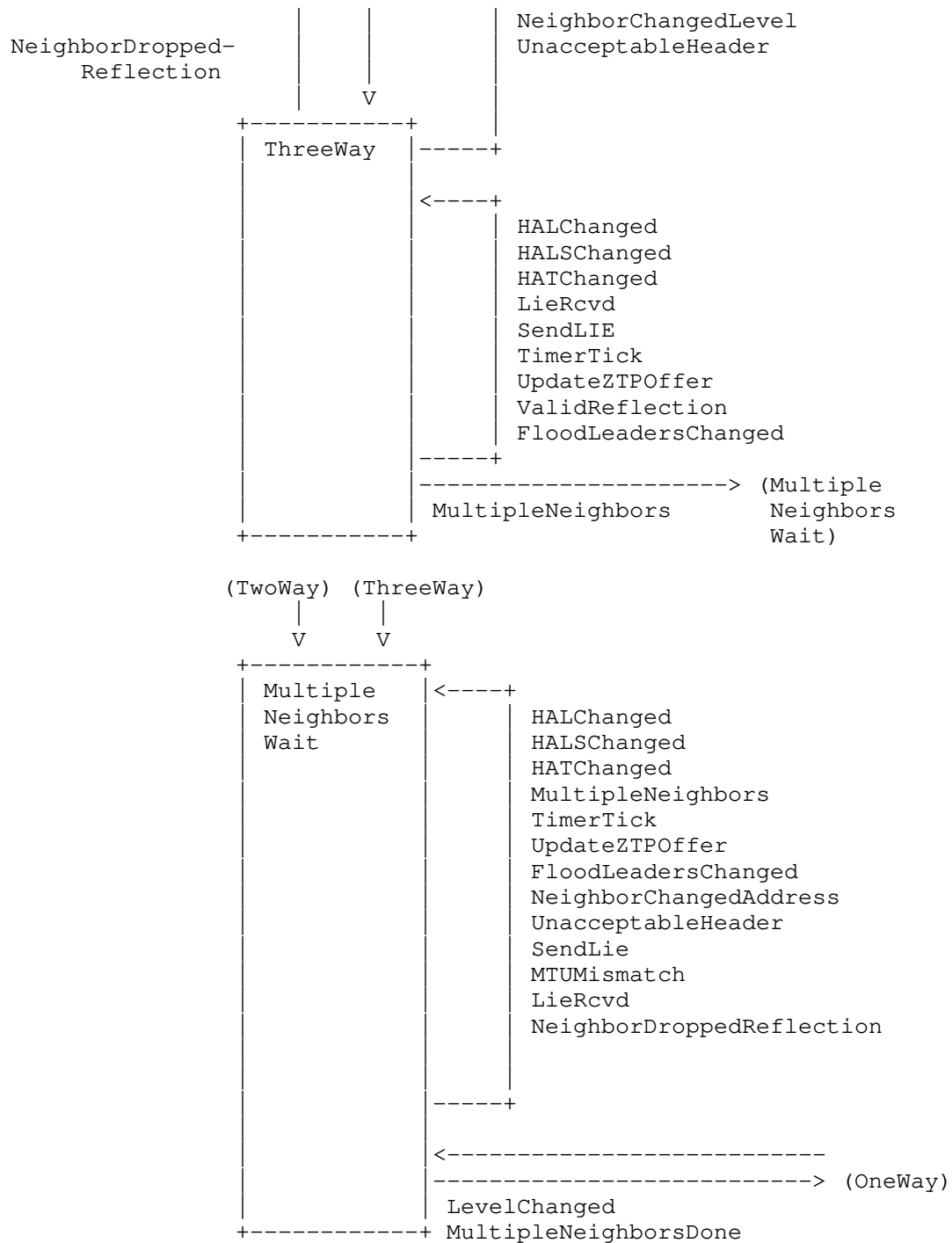


Figure 14: LIE FSM

The following words are used for well known procedures:

- * PUSH Event: queues an event to be executed by the FSM upon exit of this action
- * CLEANUP: neighbor MUST be reset to unknown
- * SEND_LIE: create and send 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_not_flood_repeater' to computed value
- * PROCESS_LIE:
 1. if LIE has major version not equal to this node's *or* system ID equal to this node's system ID or 'IllegalSystemID' then CLEANUP else
 2. if LIE has non matching MTUs then CLEANUP, PUSH UpdateZTPOffer, PUSH MTUMismatch else
 3. if LIE has undefined level OR this node's level is undefined OR this node is a leaf and remote level is lower than HAT OR (LIE's level is not leaf AND its difference is more than one from this node's level) then CLEANUP, PUSH UpdateZTPOffer, PUSH UnacceptableHeader else
 4. 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
 5. CHECK_THREE_WAY
- * CHECK_THREE_WAY: if current state is OneWay do nothing else
 1. if LIE packet does not contain neighbor then if current state is ThreeWay then PUSH NeighborDroppedReflection else
 2. if packet reflects this system's ID and local port and state is ThreeWay then PUSH event ValidReflection else PUSH event MultipleNeighbors

States:

- * OneWay: initial state FSM is starting from. In this state the neighbors did not see any valid LIEs from a neighbor after the state was entered.
- * TwoWay: that state is entered when a node has seen a LIE from a neighbor but it did not contain its reflection.
- * ThreeWay: this state signifies that lies from a neighbor are seen with correct reflection. On achieving this state the link can be advertised in 'neighbors' element in 'NodeTIEElement'.
- * MultipleNeighborsWait: 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.

Events:

- * TimerTick: one second timer tic, i.e. the event is generated for FSM by some external entity once a second. To be quietly ignored if transition does not exist.
- * LevelChanged: node's level has been changed by ZTP or configuration. This is provided by the ZTP FSM.
- * HALChanged: best HAL computed by ZTP has changed. This is provided by the ZTP FSM.
- * HATChanged: HAT computed by ZTP has changed. This is provided by the ZTP FSM.

- * **HALSChanged:** set of HAL offering systems computed by ZTP has changed. This is provided by the ZTP FSM.
- * **LieRcvd:** received LIE on the interface.
- * **NewNeighbor:** new neighbor seen on the received LIE.
- * **ValidReflection:** received reflection of this node from neighbor, i.e. 'neighbor' element in 'LiePacket' corresponds to this node.
- * **NeighborDroppedReflection:** lost previously seen reflection from neighbor, i.e. 'neighbor' element in 'LiePacket' does not correspond to this node or is not present.
- * **NeighborChangedLevel:** neighbor changed advertised level from the previously seen one.
- * **NeighborChangedAddress:** neighbor changed IP address, i.e. LIE has been received from an address different from previous LIEs. Those changes will influence the sockets used to listen to TIEs, TIREs, TIDEs.
- * **UnacceptableHeader:** Unacceptable header seen.
- * **MTUMismatch:** MTU mismatched.
- * **NeighborChangedMinorFields:** minor fields changed in neighbor's LIE.
- * **HoldtimeExpired:** adjacency holddown timer expired.
- * **MultipleNeighbors:** more than one neighbor seen on interface
- * **MultipleNeighborsDone:** multiple neighbors timer expired.
- * **FloodLeadersChanged:** node's election algorithm determined new set of flood leaders.
- * **SendLie:** send a LIE out.
- * **UpdateZTPOffer:** update this node's ZTP offer. This is sent to the ZTP FSM.

Actions:

- * on **TimerTick** in **OneWay** finishes in **OneWay**: PUSH **SendLie** event
- * on **UnacceptableHeader** in **OneWay** finishes in **OneWay**: no action

- * on LevelChanged in OneWay finishes in OneWay: update level with event value, PUSH SendLie event
- * on NeighborChangedMinorFields in OneWay finishes in OneWay: no action
- * on NeighborChangedLevel in OneWay finishes in OneWay: no action
- * on NewNeighbor in OneWay finishes in TwoWay: PUSH SendLie event
- * on HoldtimeExpired in OneWay finishes in OneWay: no action
- * on HALSChanged in OneWay finishes in OneWay: store HALS
- * on NeighborChangedAddress in OneWay finishes in OneWay: no action
- * on LieRcvd in OneWay finishes in OneWay: PROCESS_LIE
- * on ValidReflection in OneWay finishes in ThreeWay: no action
- * on SendLie in OneWay finishes in OneWay: SEND_LIE
- * on UpdateZTPOffer in OneWay finishes in OneWay: send offer to ZTP FSM
- * on HATChanged in OneWay finishes in OneWay: store HAT
- * on MultipleNeighbors in OneWay finishes in MultipleNeighborsWait: start multiple neighbors timer with interval 'multiple_neighbors_lie_holdtime_multipler' * 'default_lie_holdtime'
- * on MTUMismatch in OneWay finishes in OneWay: no action
- * on FloodLeadersChanged in OneWay finishes in OneWay: update 'you_are_flood_repeater' LIE elements based on flood leader election results
- * on NeighborDroppedReflection in OneWay finishes in OneWay: no action
- * on HALChanged in OneWay finishes in OneWay: store new HAL
- * on NeighborChangedAddress in TwoWay finishes in OneWay: no action
- * on LieRcvd in TwoWay finishes in TwoWay: PROCESS_LIE

- * on UpdateZTPOffer in TwoWay finishes in TwoWay: send offer to ZTP FSM
- * on HoldtimeExpired in TwoWay finishes in OneWay: no action
- * on MTUMismatch in TwoWay finishes in OneWay: no action
- * on UnacceptableHeader in TwoWay finishes in OneWay: no action
- * on ValidReflection in TwoWay finishes in ThreeWay: no action
- * on SendLie in TwoWay finishes in TwoWay: SEND_LIE
- * on HATChanged in TwoWay finishes in TwoWay: store HAT
- * on HALChanged in TwoWay finishes in TwoWay: store new HAL
- * on LevelChanged in TwoWay finishes in TwoWay: update level with event value
- * on FloodLeadersChanged in TwoWay finishes in TwoWay: update 'you_are_flood_repeater' LIE elements based on flood leader election results
- * on NewNeighbor in TwoWay finishes in MultipleNeighborsWait: PUSH SendLie event
- * on TimerTick in TwoWay finishes in TwoWay: PUSH SendLie event, if last valid LIE was received more than 'holdtime' ago as advertised by neighbor then PUSH HoldtimeExpired event
- * on NeighborChangedLevel in TwoWay finishes in OneWay: no action
- * on MultipleNeighbors in TwoWay finishes in MultipleNeighborsWait: start multiple neighbors timer with interval 'multiple_neighbors_lie_holdtime_multipler' * 'default_lie_holdtime'
- * on HALSChanged in TwoWay finishes in TwoWay: store HALS
- * on NeighborChangedAddress in ThreeWay finishes in OneWay: no action
- * on ValidReflection in ThreeWay finishes in ThreeWay: no action
- * on HoldtimeExpired in ThreeWay finishes in OneWay: no action
- * on UnacceptableHeader in ThreeWay finishes in OneWay: no action

- * on NeighborDroppedReflection in ThreeWay finishes in TwoWay: no action
- * on HALChanged in ThreeWay finishes in ThreeWay: store new HAL
- * on MultipleNeighbors in ThreeWay finishes in MultipleNeighborsWait: start multiple neighbors timer with interval 'multiple_neighbors_lie_holdtime_multiplier' * 'default_lie_holdtime'
- * on LevelChanged in ThreeWay finishes in OneWay: update level with event value
- * on HALSChanged in ThreeWay finishes in ThreeWay: store HALS
- * on TimerTick in ThreeWay finishes in ThreeWay: PUSH SendLie event, if last valid LIE was received more than 'holdtime' ago as advertised by neighbor then PUSH HoldtimeExpired event
- * on HATChanged in ThreeWay finishes in ThreeWay: store HAT
- * on UpdateZTPOffer in ThreeWay finishes in ThreeWay: send offer to ZTP FSM
- * on LieRcvd in ThreeWay finishes in ThreeWay: PROCESS_LIE
- * on NeighborChangedLevel in ThreeWay finishes in OneWay: no action
- * on SendLie in ThreeWay finishes in ThreeWay: SEND_LIE
- * on FloodLeadersChanged in ThreeWay finishes in ThreeWay: update 'you_are_flood_repeater' LIE elements based on flood leader election results, PUSH SendLie
- * on MTUMismatch in ThreeWay finishes in OneWay: no action
- * on HoldtimeExpired in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on LieRcvd in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on NeighborDroppedReflection in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on MTUMismatch in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action

- * on NeighborChangedBFDCapability in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on LevelChanged in MultipleNeighborsWait finishes in OneWay: update level with event value
- * on SendLie in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on UpdateZTPOffer in MultipleNeighborsWait finishes in MultipleNeighborsWait: send offer to ZTP FSM
- * on MultipleNeighborsDone in MultipleNeighborsWait finishes in OneWay: no action
- * on HATChanged in MultipleNeighborsWait finishes in MultipleNeighborsWait: store HAT
- * on NeighborChangedAddress in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on HALSChanged in MultipleNeighborsWait finishes in MultipleNeighborsWait: store HALS
- * on HALChanged in MultipleNeighborsWait finishes in MultipleNeighborsWait: store new HAL
- * on MultipleNeighbors in MultipleNeighborsWait finishes in MultipleNeighborsWait: start multiple neighbors timer with interval 'multiple_neighbors_lie_holdtime_multiplier' * 'default_lie_holdtime'
- * on FloodLeadersChanged in MultipleNeighborsWait finishes in MultipleNeighborsWait: update 'you_are_flood_repeater' LIE elements based on flood leader election results
- * on ValidReflection in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on TimerTick in MultipleNeighborsWait finishes in MultipleNeighborsWait: check MultipleNeighbors timer, if timer expired PUSH MultipleNeighborsDone
- * on UnacceptableHeader in MultipleNeighborsWait finishes in MultipleNeighborsWait: no action
- * on Entry into OneWay: CLEANUP

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.

The TIE exchange mechanism uses the port indicated by each node in the LIE exchange as `'flood_port'` in `'LIEPacket'` and the interface on which the adjacency has been formed as destination. It SHOULD use TTL of 1 or 255 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 an implementation 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 starting with `'TIEPacket'` root.

4.2.3.2. Southbound and Northbound TIE 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/east-west. Those differing TIE databases are called 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. Section 4.2.5 explains further details.

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. The direction itself is carried in `'direction'` of `'TIEID'` schema element.

| 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, 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) is considered. Unnumbered interfaces are implicitly assumed and for simplicity, the key value elements which may be included in their South TIEs or North TIEs are not shown. First, in Figure 15 are the TIEs generated by some nodes.

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 15: 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 further elaborated on in Section 4.2.3.9 and Section 4.3.7.

For node TIEs to carry more adjacencies than fit into an MTU, the element 'neighbors' may contain different set of neighbors in each TIE. Those disjoint sets of neighbors MUST be joined during according computation. Nevertheless, in case across multiple node TIEs

1. 'capabilities' do not match *or*
2. 'flags' values do not match *or*
3. same neighbor repeats with different values

the behavior is undefined and a warning SHOULD be generated after a period of time. The element 'miscabled_links' SHOULD be repeated in every node TIE, otherwise the behavior is undefined.

Different TIE types are carried in 'TIEElement'. 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 in the TIE an error MUST be reported and the TIE MUST be ignored. 'positive_disaggregation_prefixes' and 'positive_external_disaggregation_prefixes' MUST be advertised southbound only and ignored in North TIEs. 'negative_disaggregation_prefixes' MUST be aggregated and propagated according to Section 4.2.5.2 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 and ignored otherwise.

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 in link-state protocols. Although flooding is initially more demanding to implement it avoids many problems with update style used in diffused computation by distance vector protocols. However, since flooding tends to present a significant burden in large, densely meshed topologies (Fat Trees being unfortunately such a topology) RIFT provides as solution a 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 other link-state routing protocols and are largely implementation dependent.

TIEs are uniquely identified by 'TIEID' schema element. 'TIEID' space is a total order achieved by comparing the elements in sequence defined in the element and comparing each value as an unsigned integer of according length. They contain a 'seq_nr' element to distinguish newer versions of same TIE. TIEIDs also carry 'origination_time' and 'origination_lifetime'. Field 'origination_time' contains the absolute timestamp when the TIE was generated. Field 'origination_lifetime' carries lifetime when the TIE was generated. Those are normally disregarded during comparison and carried purely for debugging/security purposes if present. They may be used for comparison of last resort to differentiate otherwise equal ties and they can be used on fabrics with synchronized clock to prevent lifetime modification attacks.

Remaining lifetime counts down to 0 from origination lifetime. TIEs with lifetimes differing by less than 'lifetime_diff2ignore' MUST be considered EQUAL (if all other fields are equal). This constant MUST be larger than 'purge_lifetime' to avoid retransmissions.

All valid TIE types are defined in 'TIETypeType'. 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 with some restrictions.

4.2.3.3.1. Normative Flooding Procedures

On reception of a TIE with an undefined level value in the packet header the node MAY issue a warning and indiscriminately discard the packet. Such packets can be useful however to establish e.g. via 'instance_name', 'name' and 'originator' elements in 'LIEPacket' whether the cabling of the node fulfills expectations, even before ZTP procedures determine levels across the topology.

This section specifies the precise, normative flooding mechanism and can be omitted unless the reader is pursuing an implementation of the protocol or looks for a deep understanding of underlying information distribution mechanism.

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** either implement the given procedures in a verbatim manner or 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).

To help with adjustment of flooding speeds the encoded packets provide hints to react accordingly to losses or overruns via 'you_are_sending_too_quickly' in 'LIEPacket' and 'Packet Number' in security envelope described in Section 4.4.3. 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 packet elements and 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 "ThreeWay" state. No further tightening of this rule as to e.g. sequence is possible due to possible link buffering and re-ordering of LIEs and TIEs/TIDEs/TIREs in a real implementation for e.g. performance purposes.

A node **MUST** drop any received TIEs/TIDEs/TIREs unless it is in ThreeWay 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 on each adjacency the following elements. The word collection or queue indicates a set of elements that can be iterated over:

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 elements and procedures operating on this structure:

TIE:

Describes either a full RIFT TIE or accordingly just the 'TIEHeader' or 'TIEID' equivalent as defined in Appendix B.3. The according meaning is unambiguously contained in the context of each 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 is found on TIES_ACK then
 - a. if TIE" is 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 the following priorities if the system cannot process all the collections in real time:

1. Elements on TIES_ACK should be processed with highest priority
2. TIES_TX
3. 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 as explained in Appendix A 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 from other queues 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 TIE DB 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 computed per interface 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 sent 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 <= TIDE.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 'remaining_lifetime' field in 'TIEHeaderWithLifeTime' 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. 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.

4.2.3.3.1.6. TIEs Processing In LSDB

The Link State Database can be considered to be a switchboard that does not need any flooding procedures but can be given versions of TIEs by peers. Consecutively, after version tie-breaking by LSDB, a peer receives from the LSDB newest versions of TIEs received by other peers and processes them (without any filtering) just like receiving TIEs from its remote peer. Such a publisher model can be implemented in many ways, either in a single thread of execution or in parallel threads.

LSDB can be logically considered as the entity aging out TIEs, i.e. being responsible to discard TIEs that are stored longer than `'remaining_lifetime'` on their reception.

LSDB is also expected to periodically re-originate the node's own TIEs. It is recommended to originate at interval significantly shorter than `'default_lifetime'` to prevent TIE expiration by other nodes in the network which can lead to instabilities.

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 except 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 as described in 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 15. The flooding from particular nodes of the TIEs is given in Table 4.

| Local Node | Neighbor Node | TIEs Flooded from Local to Neighbor Node |
|------------|---------------|---|
| 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 'you_are_sending_too_quickly' flag to indicate to the neighbor that it should temporarily flood node TIEs only to it and slow down the flooding of any other TIEs. 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 indication if northbound.

Obviously this mechanism is most useful in the 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 permanently will lead to blackholes.

4.2.3.6. Initial and Periodic Database Synchronization

The initial exchange of RIFT includes periodic TIDE exchanges that contain description of the link state database and TIRES which perform the function of requesting unknown TIEs as well as confirming reception of flooded TIEs. 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, i.e. higher sequence number, empty version of the TIE with a short lifetime given by 'purge_lifetime' constant 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 generation of the content, 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 TIE content and use that as part of the 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 4.3.7.1.

A node X that

1. is **not** overloaded **and**
2. has southbound or East-West adjacencies

SHOULD originate in its south prefix TIE such a default route if and only if

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 SHOULD install a default discard route if it did not compute a default route during N-SPF. This makes the top of the fabric basically a blackhole for unreachable addresses.

4.2.3.9. Northbound TIE Flooding Reduction

RIFT chooses only a subset of northbound nodes to propagate flooding and with that both balances it (to prevent 'hot' flooding links) across the fabric as well as reduces its volume. The solution is based on several principles:

1. a node MUST flood self-originated North TIEs to all the reachable nodes at the level above which is called 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. The element 'you_are_flood_repeater' MUST be ignored if received from a northbound adjacency.

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

- * let $|NA(Node)$ be the set of Northbound adjacencies of node Node and $CN(Node)$ be the cardinality of $|NA(Node)$;
- * let $|SA(Node)$ be the set of Southbound adjacencies of node Node and $CS(Node)$ be the cardinality of $|SA(Node)$;
- * let $|P(Node)$ be the set of node Node's parents;
- * let $|G(Node)$ be the set of node Node's grandparents. Observe that $|G(Node) = |P(|P(Node))$;
- * let N be the child node at level L computing a set of FR;
- * let P be a node at level L+1 and a parent node of N, i.e. bi-directionally reachable over adjacency $ADJ(N, P)$;
- * 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 $ADJ(P, G)$;
- * let R be a redundancy constant integer; a value of 2 or higher for R is RECOMMENDED;
- * 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. $|a-b| \leq S$.
- * 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. The election can use the LIE FSM 'FloodLeadersChanged' event to notify LIE FSMs of necessity to update the sent LIEs.
3. A node MUST always flood its self-originated TIEs to all its neighbors.
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 Section 4.2.3.9, Paragraph 10, Item 6.

5. The indication of flood reduction capability MUST be carried in the node TIEs in the 'flood_reduction' element 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 where they ultimately belong. Since flooding can begin before ZTP is "finished" and in fact must do so given there is no global termination criteria for the unsynchronized ZTP algorithm, information may end up temporarily 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 the 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 TIDEs coming from the north. Those headers are then propagated southbound towards the leaf nudging it to originate a higher sequence number of the TIE effectively refreshing it all the way up to ToF.

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. N-SPF and S-SPF notation denotes here 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. This specification however uses simple, familiar SPF algorithms and concepts as example due to their prevalence in today's routing.

For reachability computation purposes RIFT considers all parallel links between two nodes to be of the same cost advertised in 'cost' element of 'NodeNeighborsTIEElement'. In case the neighbor has multiple parallel links at different cost, the largest distance (highest numerical value) MUST be advertised. Given the range of thrift encodings, 'infinite_distance' is defined as largest non-negative 'MetricType'. Any link with metric larger than that (i.e. negative MetricType) MUST be ignored in computations. Any link with metric set to 'invalid_distance' MUST be ignored in computation as well. In case of a negatively distributed prefix the metric attribute MUST be set to 'infinite_distance' by the originator and it MUST be ignored by all nodes during computation except for the purpose of determining transitive propagation and building the according routing table.

A prefix can carry the 'directly_attached' attribute to indicate that the prefix is directly attached, i.e. should be routed to even if the node is in overload. In case of a negatively distributed prefix this attribute MUST not be included by the originator and it MUST be ignored by all nodes during SPF computation. If a prefix is locally originated the attribute 'from_link' can indicate the interface to which the address belongs to. In case of a negatively distributed prefix this attribute MUST NOT be included by the originator and it MUST be ignored by all nodes during computation. A prefix can also carry the 'loopback' attribute to indicate the said property.

Prefixes are carried in different type of TIEs indicating their type. For same prefix being included in different TIE types according to Section 4.3.1. In case the same prefix is included multiple times in multiple TIEs of same type originating at the same node the resulting behavior is unspecified.

4.2.4.1. Northbound Reachability 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, the computation uses the next higher level's node South TIEs to find according adjacencies to verify backlink connectivity. Two unidirectional links MUST be associated together to confirm bidirectional connectivity, a process often known as 'backlink check'. As part of the check, both node TIEs MUST contain the correct system IDs *and* expected levels.

Default route found when crossing an E-W link SHOULD be used if and only if

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 if and only if

1. no north neighbors are advertising same or superssuming non-default prefix *and*
2. the node does not originate a non-default superssuming 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 disaggregation in pathological fabric partitioning scenarios.

A detailed example can be found in Section 5.4.

4.2.4.2. Southbound Reachability SPF

S-SPF MUST use the southbound adjacencies in the node South TIEs exclusively, i.e. progresses towards nodes at lower levels. Observe that E-W adjacencies are NEVER used in this 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 System ID 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 and when used for forwarding purposes, they defeat statistical non-blocking guarantees that Clos is providing normally.

4.2.5. Automatic Disaggregation on Link & Node Failures

4.2.5.1. Positive, Non-transitive Disaggregation

Under normal circumstances, a 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.

This specification refers to the process of advertising additional prefixes southbound as 'positive disaggregation'. Such disaggregation 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 disaggregation necessary to prevent looping or bow-tying the fabric.

A node determines the set of prefixes needing disaggregation 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 $|R$, its set of next-hops is defined 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 $|N$, its set of southbound adjacencies is defined to be $|A(n)$.
3. For a given r , if the intersection of $|H(r)$ and $|A(n)$, for any n , is empty then that prefix r must be explicitly advertised by the node in an South TIE.
4. Identical set of disaggregated 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.

Finally, a possible algorithm is described here:

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 16: 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 4.3.2 and carried in 'overload' schema element have to be respected during the computation, i.e. node advertising themselves as overloaded MUST NOT be transited in reachability computation but MUST be used as terminal nodes with prefixes they advertise being reachable.

3. all the lower level nodes are flooded the same disaggregated prefixes since RIFT does not build an South TIE per node which would 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. Different implementation strategies can be used to lessen that effect but those are clearly outside the scope of this specification.

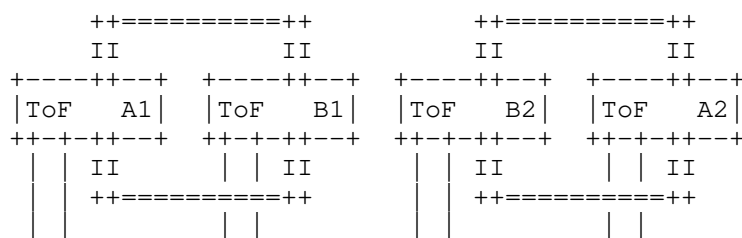
To close this section it is worth to observe that in a single plane ToF this disaggregation prevents blackholing up to $(K_LEAF * P)$ link failures in terms of Section 4.1.2 or in other terms, it takes at minimum that many link failures to partition the ToF into multiple planes.

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

Returning in this section to designs with multiple planes as shown originally in Figure 3, Figure 17 highlights now 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.



~~~ Highlighted ToF of the previous multi-plane figure ~~~

Figure 17: Topologically Connected Planes

Section 4.1.3 already describes how 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 and next section explains the involved mechanisms in more detail.

#### 4.2.5.2.2. Transitive Advertisement of Negative Disaggregates

A ToF node discovering that it cannot reach a fallen leaf SHOULD disaggregate 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. This will be described 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 horizontal 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 the southbound SPF.

The second mechanism hinges on the understanding how the negative prefixes are used within the computation as described in Figure 18. 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 4.3.4. 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 18: 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 ToP node T1 with 4 ToF parents S1..S4 as represented in Figure 19 are considered further:

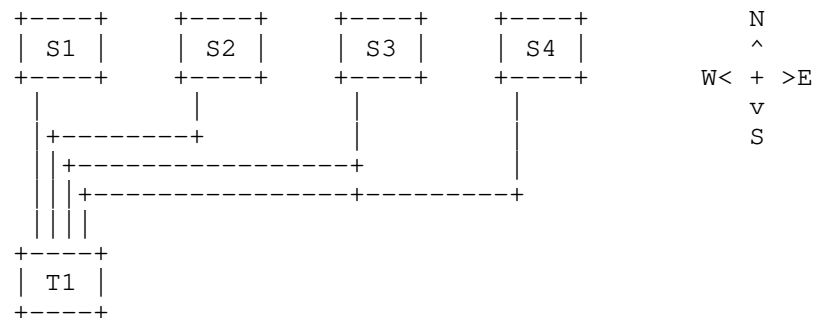


Figure 19: A ToP Node with 4 Parents

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 20:

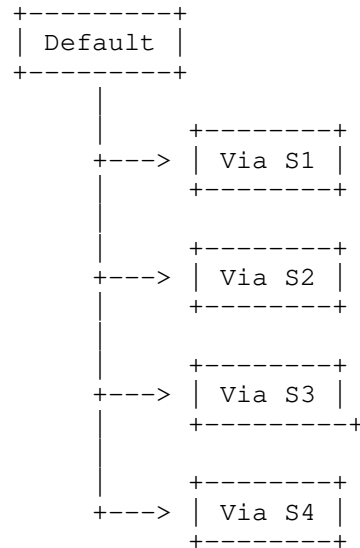


Figure 20: Abstract RIB

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 21 and Figure 22, respectively:

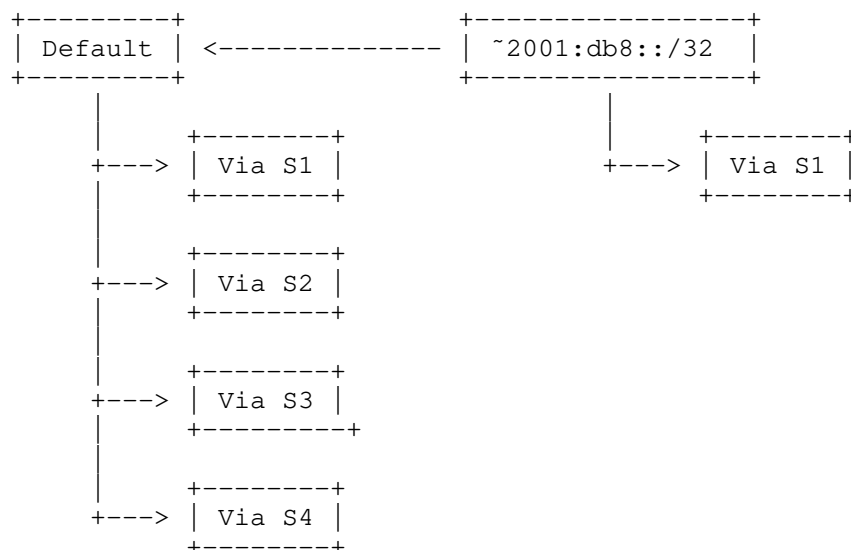


Figure 21: Abstract RIB after Negative 2001:db8::/32 from S1

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:

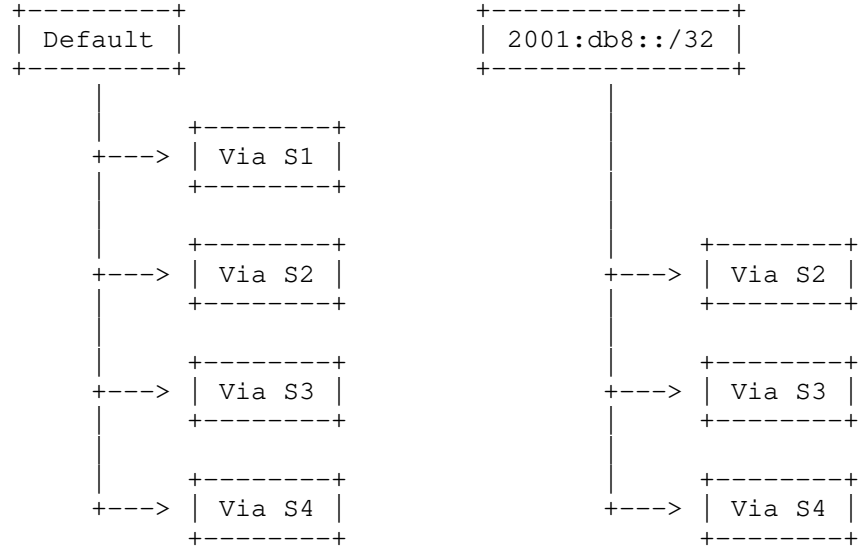


Figure 22: Abstract FIB after Negative 2001:db8::/32 from S1

To illustrate matters further 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 23:

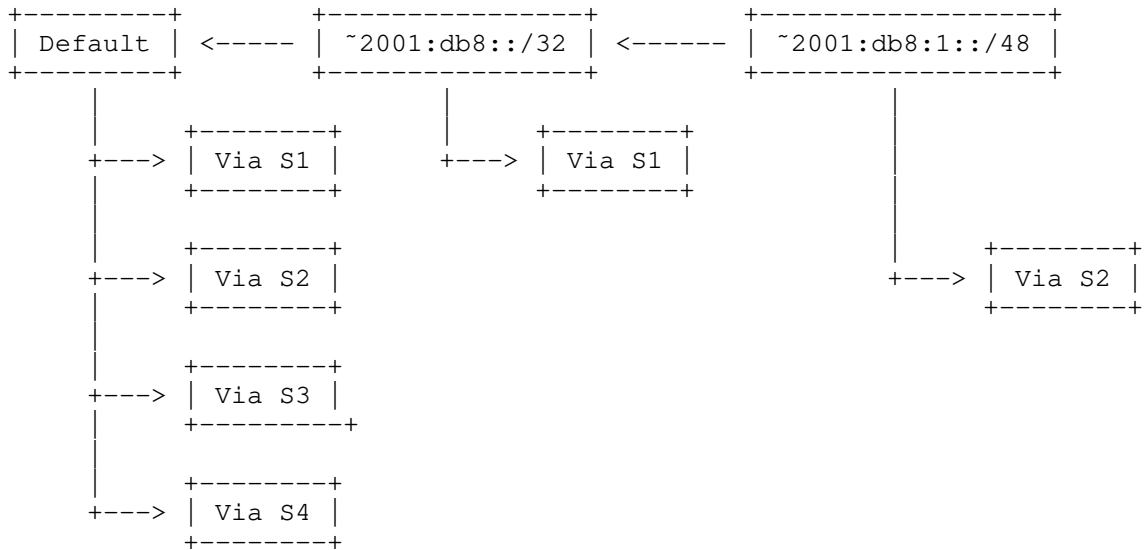


Figure 23: Abstract RIB after Negative 2001:db8:1::/48 from S2

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 24:

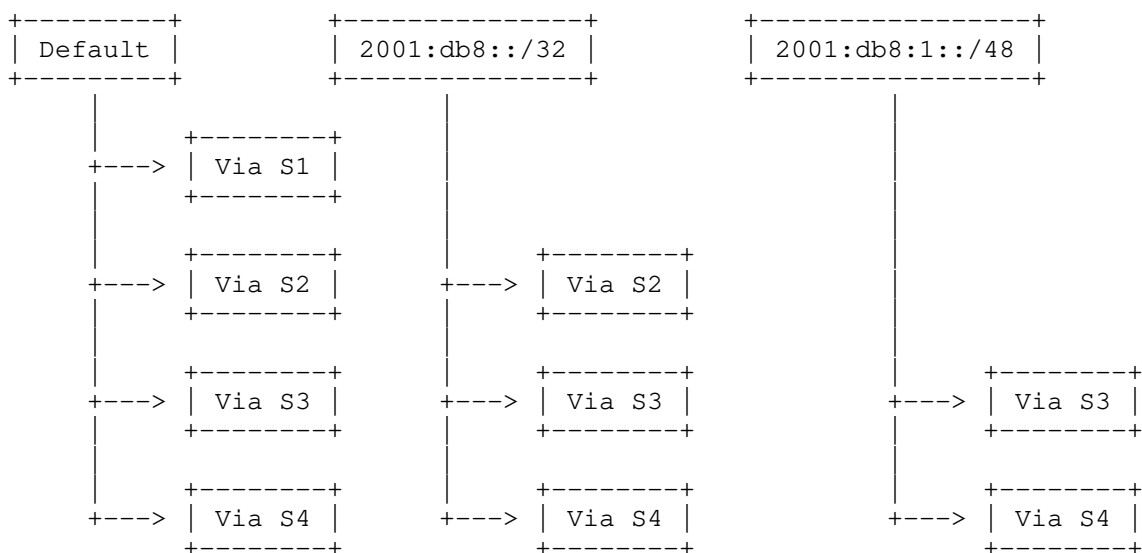


Figure 24: Abstract FIB after Negative 2001:db8:1::/48 from S2

Further, assume that S3 stops advertising its service as default gateway. The entry is removed from RIB as usual. In order to update the FIB, it is necessary to eliminate the FIB entry for the default route, as well as all the FIB entries that were created for negative routes pointing to the RIB entry being removed (::/0). This is done recursively for 2001:db8::/32 and then for, 2001:db8:1::/48. The related FIB entries via S3 are removed, as illustrated in Figure 25.

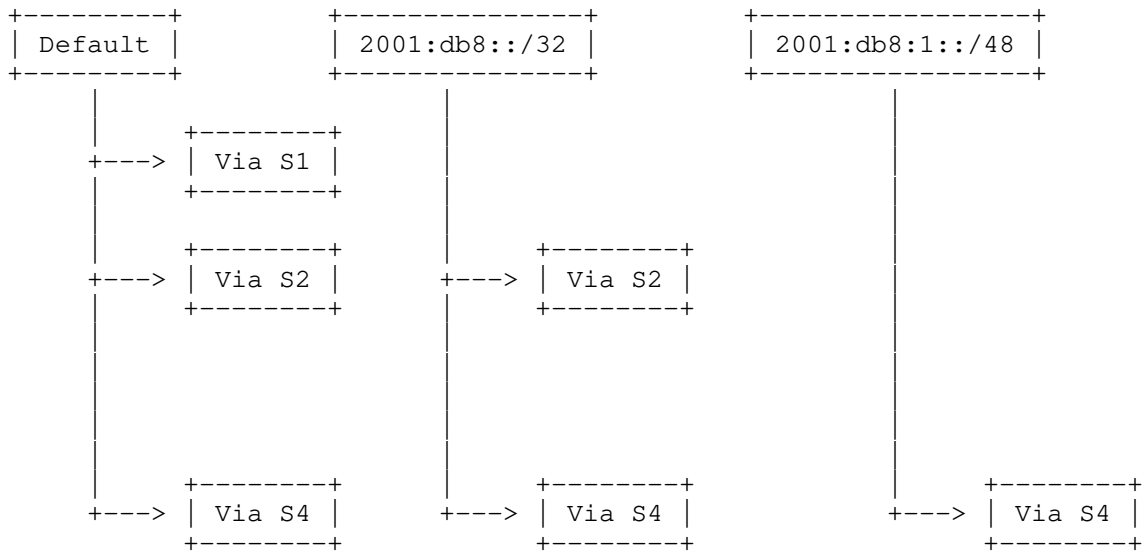


Figure 25: Abstract FIB after Loss of S3

Say that at that time, S4 would also disaggregate prefix 2001:db8:1::/48. This would mean that the FIB entry for 2001:db8:1::/48 becomes a discard route, and that would be the signal for T1 to disaggregate prefix 2001:db8:1::/48 negatively in a transitive fashion with its own children.

Finally, the case occurs where S3 becomes available again as a default gateway, and a negative advertisement is received from S4 about prefix 2001:db8:2::/48 as opposed to 2001:db8:1::/48. Again, a negative route is stored in the RIB, and the more specific route to the complementing ToF nodes are installed in FIB. Since 2001:db8:2::/48 inherits from 2001:db8::/32, the positive FIB routes are chosen by removing S4 from S2, S3, S4. The abstract FIB in T1 now shows as illustrated in Figure 26:

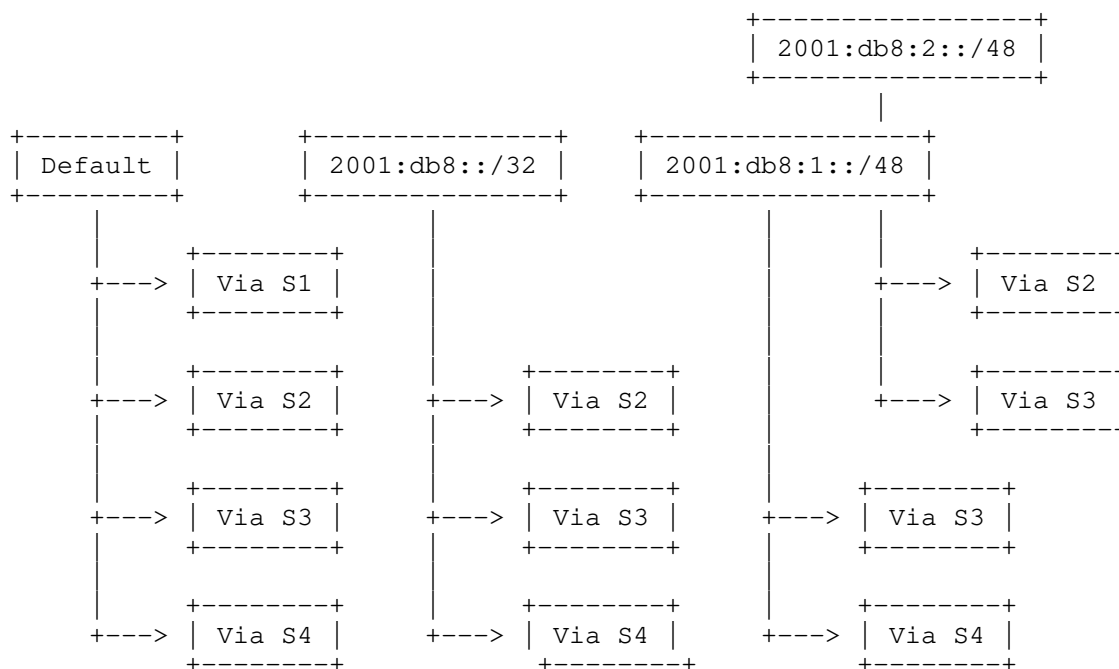


Figure 26: Abstract FIB after Negative 2001:db8:2::/48 from S4

#### 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 of offers reaches a node from "below" faster than from "above", it will greedily abandon already negotiated level derived from nodes topologically below it and properly peer with nodes above.

The fabric is very consciously numbered from the top down 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 'leaf\_level'. It is indicated in 'leaf\_only' schema element.

##### 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). It is carried in 'top\_of\_fabric' schema element. A standards conform RIFT implementation implies a CONFIGURED\_LEVEL value of 'top\_of\_fabric\_level' in case of TOP\_OF\_FABRIC. This value is kept reasonably low to allow for fast ZTP re-convergence on failures.

##### 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 in the sense of deriving its own level based on other nodes' information. 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 'leaf\_level'.

**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 4.3.9. 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. It is carried in the 'leaf\_only\_and\_leaf\_2\_leaf\_procedures' schema 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 'leaf\_level' do not constitute VOLs (since no valid DERIVED\_LEVEL can be obtained from those and consequently 'not\_a\_ztp\_offer' flag 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 ThreeWay (HAT):**

Highest neighbor level of all the formed ThreeWay adjacencies for the node.



#### 4.2.7.2. Automatic System ID Selection

RIFT nodes require a 64 bit System ID 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, i.e. purging them) 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 considerations of miscabled or unusually cabled fabric and how such a topology can be forced into a "lattice" structure which a fabric represents (with further restrictions). A necessary and sufficient physical cabling is shown in Figure 27. The assumption here is that all nodes are in the same PoD.

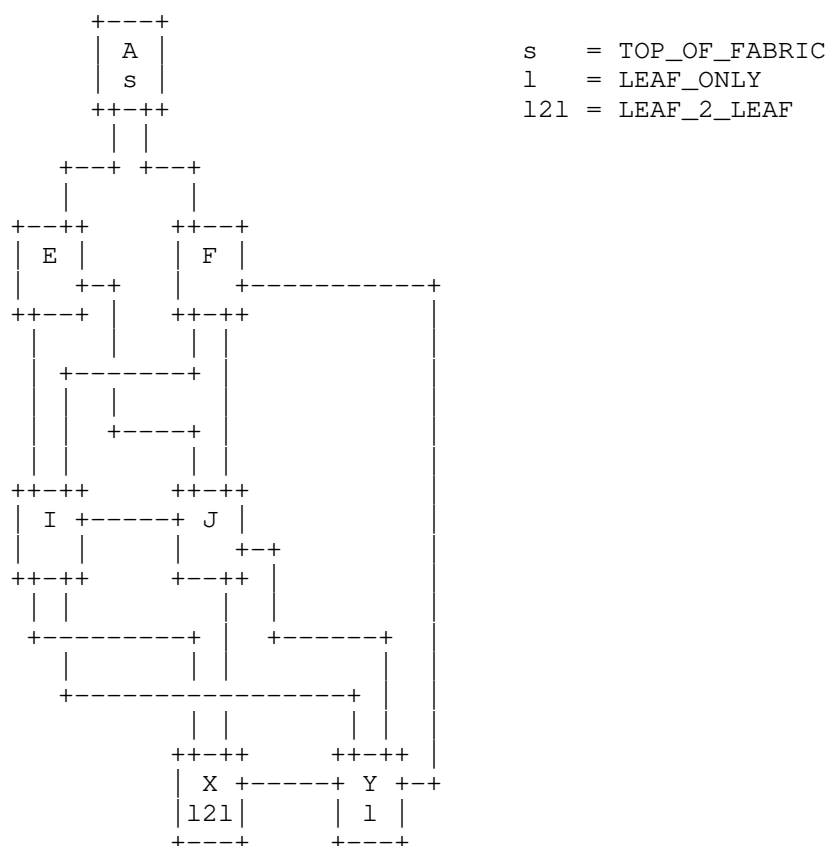


Figure 27: Generic ZTP Cabling Considerations

First, RIFT 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 the protocol has 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 the protocol chooses 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 timer 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 ThreeWay 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:

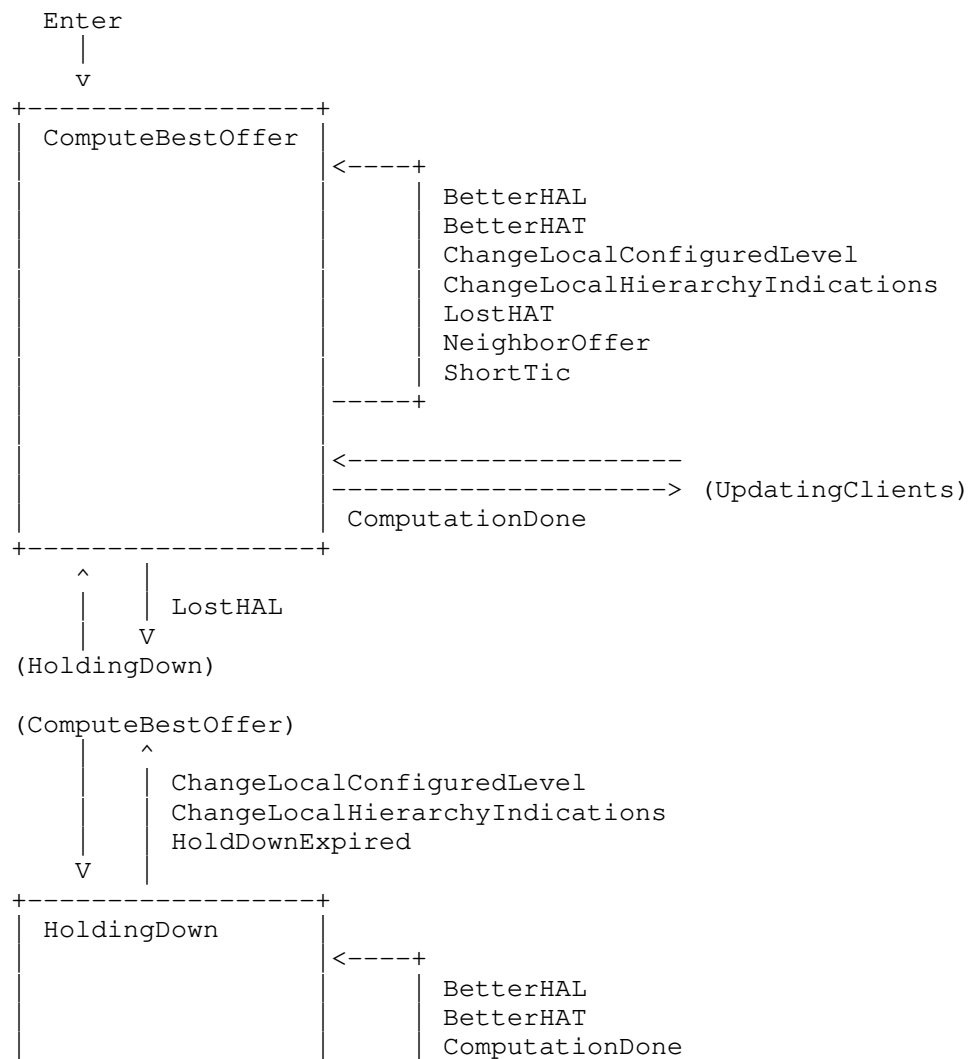
1. It may pick a different strategy to choose VOL, e.g. use the VOL value with highest number of VOLs. Such strategies are only possible since the node always remains "at the bottom of the

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.



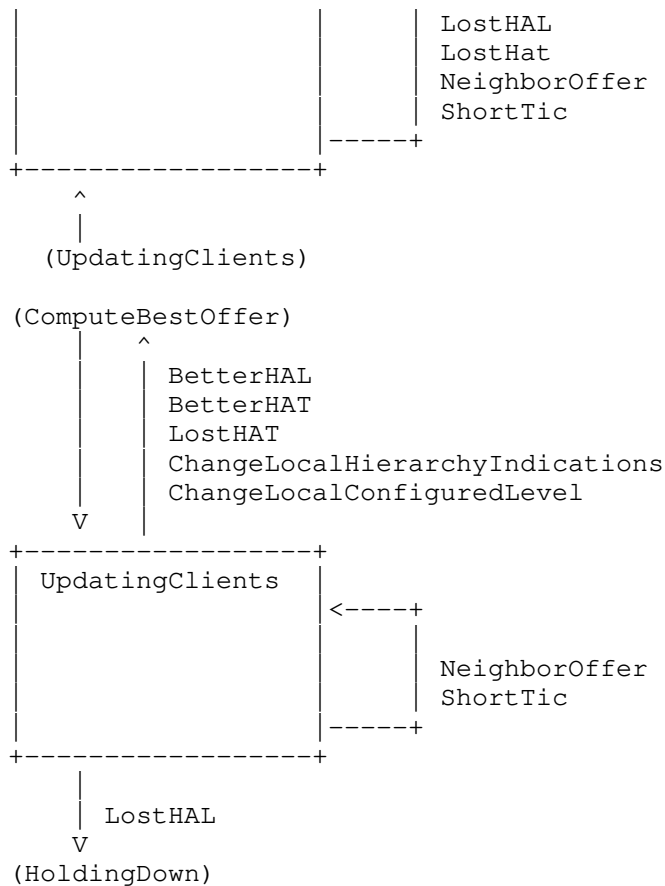


Figure 28: ZTP FSM

The following words are used for well known procedures:

- \* PUSH Event: queues an event to be executed by the FSM upon exit of this action
- \* COMPARE\_OFFERS: checks whether based on current offers and held last results the events BetterHAL/LostHAL/BetterHAT/LostHAT are necessary and returns them
- \* UPDATE\_OFFER: store current offer with adjacency holdtime as lifetime and COMPARE\_OFFERS, then PUSH according events
- \* LEVEL\_COMPUTE: compute best offered or configured level and HAL/HAT, if anything changed PUSH ComputationDone

- \* REMOVE\_OFFER: remove the according offer and COMPARE\_OFFERS, PUSH according events
- \* PURGE\_OFFERS: REMOVE\_OFFER for all held offers, COMPARE OFFERS, PUSH according events
- \* PROCESS\_OFFER:
  - 1. if no level offered then REMOVE\_OFFER
  - 2. else
    - 1. if offered level > leaf then UPDATE\_OFFER
    - 2. else REMOVE\_OFFER

States:

- \* ComputeBestOffer: processes received offers to derive ZTP variables
- \* HoldingDown: holding down while receiving updates
- \* UpdatingClients: updates other FSMs with computation results

Events:

- \* ChangeLocalHierarchyIndications: node locally configured with new leaf flags.
- \* ChangeLocalConfiguredLevel: node locally configured with a defined level
- \* NeighborOffer: a new neighbor offer with optional level and neighbor state.
- \* BetterHAL: better HAL computed internally.
- \* BetterHAT: better HAT computed internally.
- \* LostHAL: lost last HAL in computation.
- \* LostHAT: lost HAT in computation.
- \* ComputationDone: computation performed.
- \* HoldDownExpired: holddown timer expired.

- \* ShortTic: one second timer tic, i.e. the event is generated for FSM by some external entity once a second. To be ignored if transition does not exist.

Actions:

- \* on ChangeLocalConfiguredLevel in HoldingDown finishes in ComputeBestOffer: store configured level
- \* on BetterHAT in HoldingDown finishes in HoldingDown: no action
- \* on ShortTic in HoldingDown finishes in HoldingDown: remove expired offers and if holddown timer expired PUSH\_EVENT HoldDownExpired
- \* on NeighborOffer in HoldingDown finishes in HoldingDown: PROCESS\_OFFER
- \* on ComputationDone in HoldingDown finishes in HoldingDown: no action
- \* on BetterHAL in HoldingDown finishes in HoldingDown: no action
- \* on LostHAT in HoldingDown finishes in HoldingDown: no action
- \* on LostHAL in HoldingDown finishes in HoldingDown: no action
- \* on HoldDownExpired in HoldingDown finishes in ComputeBestOffer: PURGE\_OFFERS
- \* on ChangeLocalHierarchyIndications in HoldingDown finishes in ComputeBestOffer: store leaf flags
- \* on LostHAT in ComputeBestOffer finishes in ComputeBestOffer: LEVEL\_COMPUTE
- \* on NeighborOffer in ComputeBestOffer finishes in ComputeBestOffer: PROCESS\_OFFER
- \* on BetterHAT in ComputeBestOffer finishes in ComputeBestOffer: LEVEL\_COMPUTE
- \* on ChangeLocalHierarchyIndications in ComputeBestOffer finishes in ComputeBestOffer: store leaf flags and LEVEL\_COMPUTE
- \* 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 ShortTic in ComputeBestOffer finishes in ComputeBestOffer: remove expired offers
- \* on ComputationDone in ComputeBestOffer finishes in UpdatingClients: no action
- \* on ChangeLocalConfiguredLevel in ComputeBestOffer finishes in ComputeBestOffer: store configured level and LEVEL\_COMPUTE
- \* on BetterHAL in ComputeBestOffer finishes in ComputeBestOffer: LEVEL\_COMPUTE
- \* on ShortTic in UpdatingClients finishes in UpdatingClients: remove expired offers
- \* 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 BetterHAT in UpdatingClients finishes in ComputeBestOffer: no action
- \* on BetterHAL in UpdatingClients finishes in ComputeBestOffer: no action
- \* on ChangeLocalConfiguredLevel in UpdatingClients finishes in ComputeBestOffer: store configured level
- \* on ChangeLocalHierarchyIndications in UpdatingClients finishes in ComputeBestOffer: store leaf flags
- \* on NeighborOffer in UpdatingClients finishes in UpdatingClients: PROCESS\_OFFER
- \* on LostHAT in UpdatingClients finishes in ComputeBestOffer: no action
- \* on Entry into ComputeBestOffer: LEVEL\_COMPUTE
- \* on Entry into UpdatingClients: update all LIE FSMs with computation results

#### 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 29.



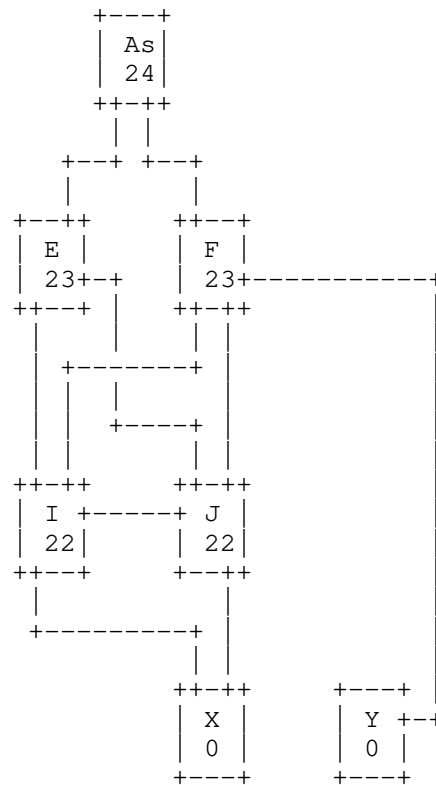


Figure 29: Generic ZTP Topology Autoconfigured

In case where the LEAF\_ONLY restriction on Y is removed the outcome would be very different however and result in Figure 30. This demonstrates basically that auto configuration makes miscabling detection hard and with that can lead to undesirable effects in cases where leaves are not "nailed" by the accordingly configured flags and arbitrarily cabled.

A node MAY analyze the outstanding level offers on its interfaces and generate warnings when its internal ruleset flags a possible miscabling. As an example, when a node's sees ZTP level offers that differ by more than one level from its chosen level (with proper accounting for leaf's being at level 'leaf\_level') this can indicate miscabling.

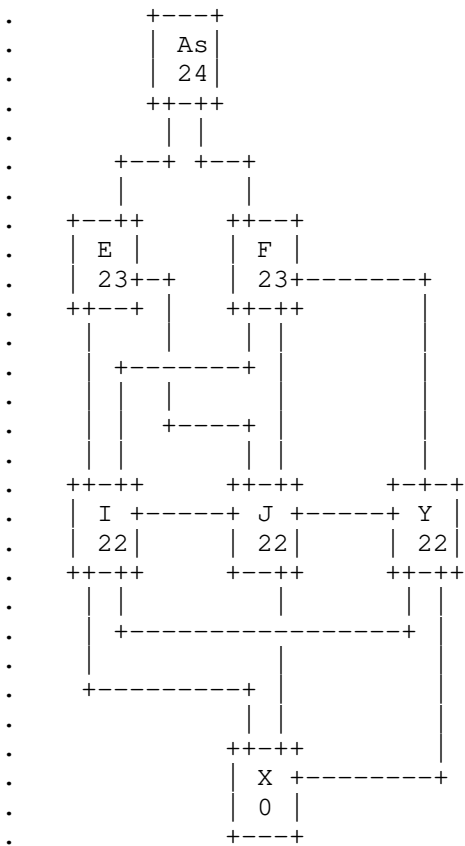


Figure 30: Generic ZTP Topology Autoconfigured

4.3. Further Mechanisms

4.3.1. Route Preferences

Since RIFT distinguishes between different route types such as e.g. external routes from other protocols and additionally advertises special types of routes on disaggregation, the protocol MUST tie-break internally different types on a clear preference scale to prevent blackholes or loops. The preferences are given in the schema 'RouteType'.

Table Table 5 contains the route type as derived from the TIE type carrying it from the most preferred to the least preferred one.

| TIE Type                                                         | Resulting Route Type |
|------------------------------------------------------------------|----------------------|
| None                                                             | Discard              |
| Local Interface                                                  | LocalPrefix          |
| S-PGP                                                            | South PGP            |
| N-PGP                                                            | North PGP            |
| North Prefix                                                     | NorthPrefix          |
| North External Prefix                                            | NorthExternalPrefix  |
| South Prefix and South Positive Disaggregation                   | SouthPrefix          |
| South External Prefix and South Positive External Disaggregation | SouthExternalPrefix  |
| South Negative Prefix                                            | NegativeSouthPrefix  |

Table 5: TIEs and Contained Route Types

#### 4.3.2. Overload Bit

Overload attribute is specified in the packet encoding schema (Appendix B).

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 attribute 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.

#### 4.3.3. 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.

#### 4.3.4. 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 is also called a 'monotonic clock' in the schema) that consists of a timestamp and optional sequence number field. In case of a negatively distributed prefix this attribute MUST NOT be included by the originator and it MUST be ignored by all nodes during computation. 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):

- \* 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.
- \* 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.
- \* 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.
- \* Any prefix present on the fabric in multiple nodes that has the 'same' clock is considered as anycast.
- \* 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, and provides for this purpose a variable called MAXIMUM\_CLOCK\_DELTA.

#### 4.3.4.1. 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].

#### 4.3.4.2. 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.

#### 4.3.4.3. 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 4.3.4.1.

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.

#### 4.3.4.4. 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.

#### 4.3.5. Key/Value Store

##### 4.3.5.1. 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.

#### 4.3.5.2. 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, the according mechanism will define, if necessary, according tie-breaking rules depending on the semantics of the information.

Only KV TIEs from nodes that are reachable via multiplane reachability computation mentioned in Section 4.2.5.2.3 SHOULD be considered.

#### 4.3.6. Interactions with BFD

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

After RIFT ThreeWay 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 'LinkCapabilities'.

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].

#### 4.3.7. 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.

##### 4.3.7.1. 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. The bandwidth is advertised in 'NodeNeighborsTIEElement' element which represents the sum of the bandwidths of all the parallel links to a neighbor. Default routes with differing distances SHOULD be used to support weighted ECMP forwarding. Such a distance is called Bandwidth Adjusted Distance or BAD. This is best illustrated by a simple example.

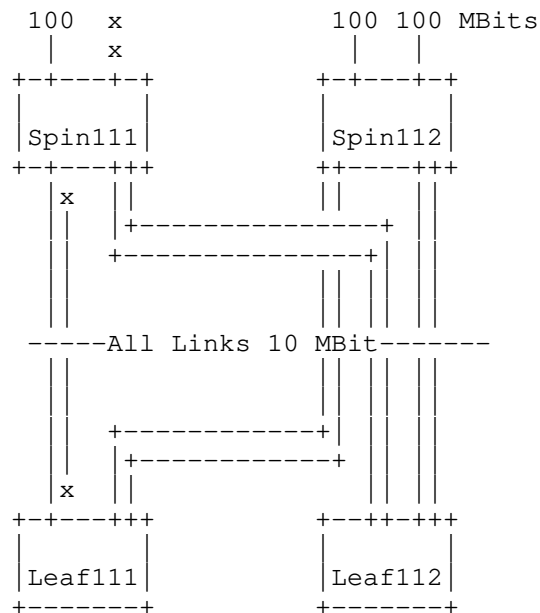


Figure 31: Balancing Bandwidth

Figure 31 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. It includes parallel link failure between Leaf 111 and Spine 111 and as a result, Leaf 111 wants to forward more traffic toward Spine 112. Additionally, it includes as well 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 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  $BAD = D * (1 + MAX\_M\_N\_u - M\_N\_u)$  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. The implicit assumption here is that all default route distances are advertised with  $D=1$  and that  $OVERSUBSCRIPTION\_CONSTANT = 1$ .

| 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 6: 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.

#### 4.3.7.2. 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.

#### 4.3.8. 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.

#### 4.3.9. 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).

#### 4.3.10. 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 4.3.6 are implementation dependent when multiple RIFT instances run on the same link.

#### 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.

- \* The most security conscious operators may choose to have control over which ports interconnect between a given pair of nodes, such a model is called 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.
- \* In physically secure data center locations, operators may choose to control connectivity between entire nodes, called here the "Node-Association Model" (NAM). A benefit of this model is that it allows for simplified port sparing.
- \* In the most relaxed environments, an operator may only choose to control which nodes join a particular fabric. This is denoted as the "Fabric-Association Model" (FAM). This is achievable by using a single shared secret across the entire fabric. Such flexibility makes sense when servers are considered as leaf devices, as those are replaced more often than network nodes. In addition, this model allows for simplified node sparing.
- \* 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 32 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.

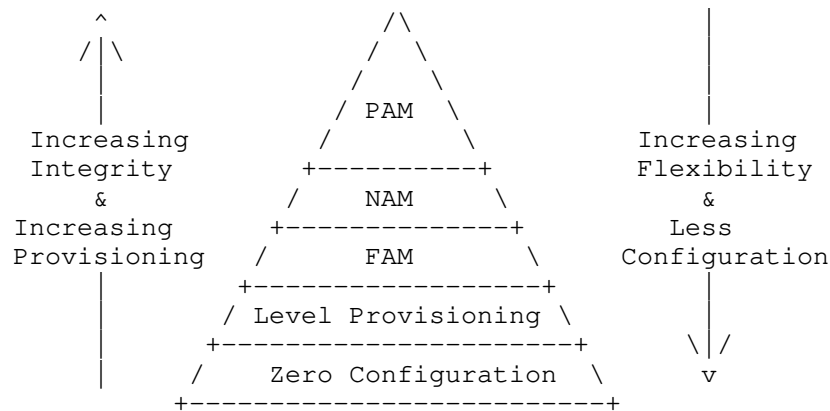


Figure 32: Security Model

#### 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

A serialized schema 'ProtocolPacket' MUST be carried in a secure envelope illustrated in Figure 33. Any value in the packet following a security fingerprint MUST be used only after the appropriate fingerprint has been validated against the data covered by it and advertised key.

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

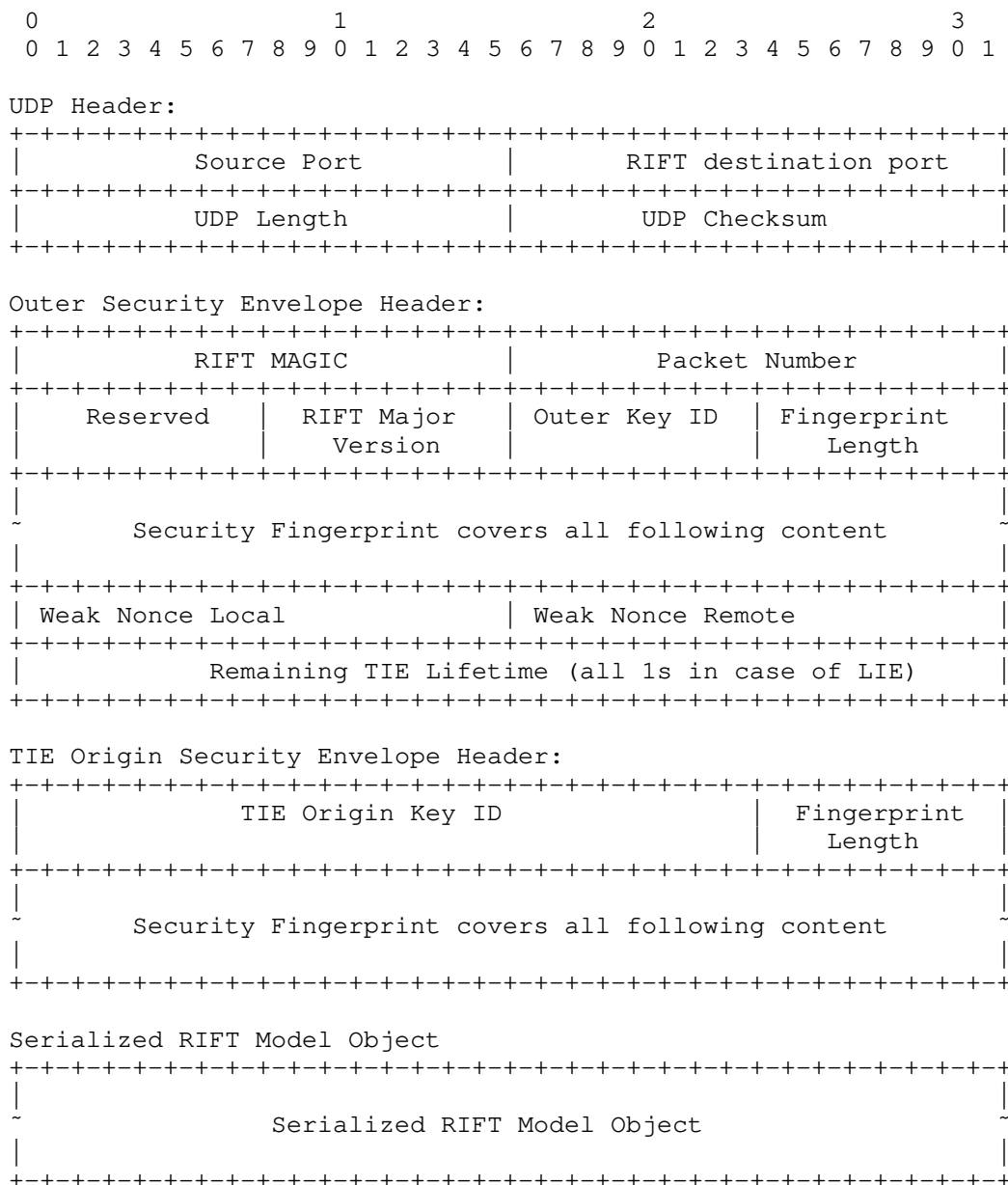


Figure 33: Security Envelope

**RIFT MAGIC:**

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



**Packet Number:**

16 bits. An optional, per adjacency, per packet type monotonically increasing 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. Packet numbers are incremented on each interface and within that for each type of packet independently. This allows to parallelize packet generation and processing for different types within an implementation if so desired.

**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 'invalid\_key\_value\_key' 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 'invalid\_key\_value\_key' means that no valid fingerprint was computed. This key ID scope is global to the RIFT instance since it may imply the originator of the TIE so the contained object does not have to be de-serialized to obtain the originator.

**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.

**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. The term "nonce" is used a bit loosely since RIFT nonces are not being changed in every packet as often common in cryptography. For efficiency purposes they are changed at a high enough frequency to dwarf practical replay attack attempts. And hence, such nonces are called from this point on "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 one and MUST increment its nonce at least every `'nonce_regeneration_interval'` (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 in `'maximum_valid_nonce_delta'` 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. The delta is either a negative or positive difference that a mirrored nonce can deviate from local value to be considered valid. If nonces are not changed on every packet but at the maximum interval on both sides this opens statistically a `'maximum_valid_nonce_delta'/2` window of identical LIEs, TIE and TI(x)E replays. The interval cannot be too small since LIE FSM may change states fairly quickly during ZTP without sending LIEs and additionally, UDP can both loose as well as misorder packets.

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 TwoWay or ThreeWay state and accept an `'undefined_nonce'` for its local nonce value on packets in any other state than ThreeWay.

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 ThreeWay to OneWay 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'. This will limit the number of computations performed for security purposes caused by lifetime attacks as long the weak nonce did not advance.

#### 4.5. 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

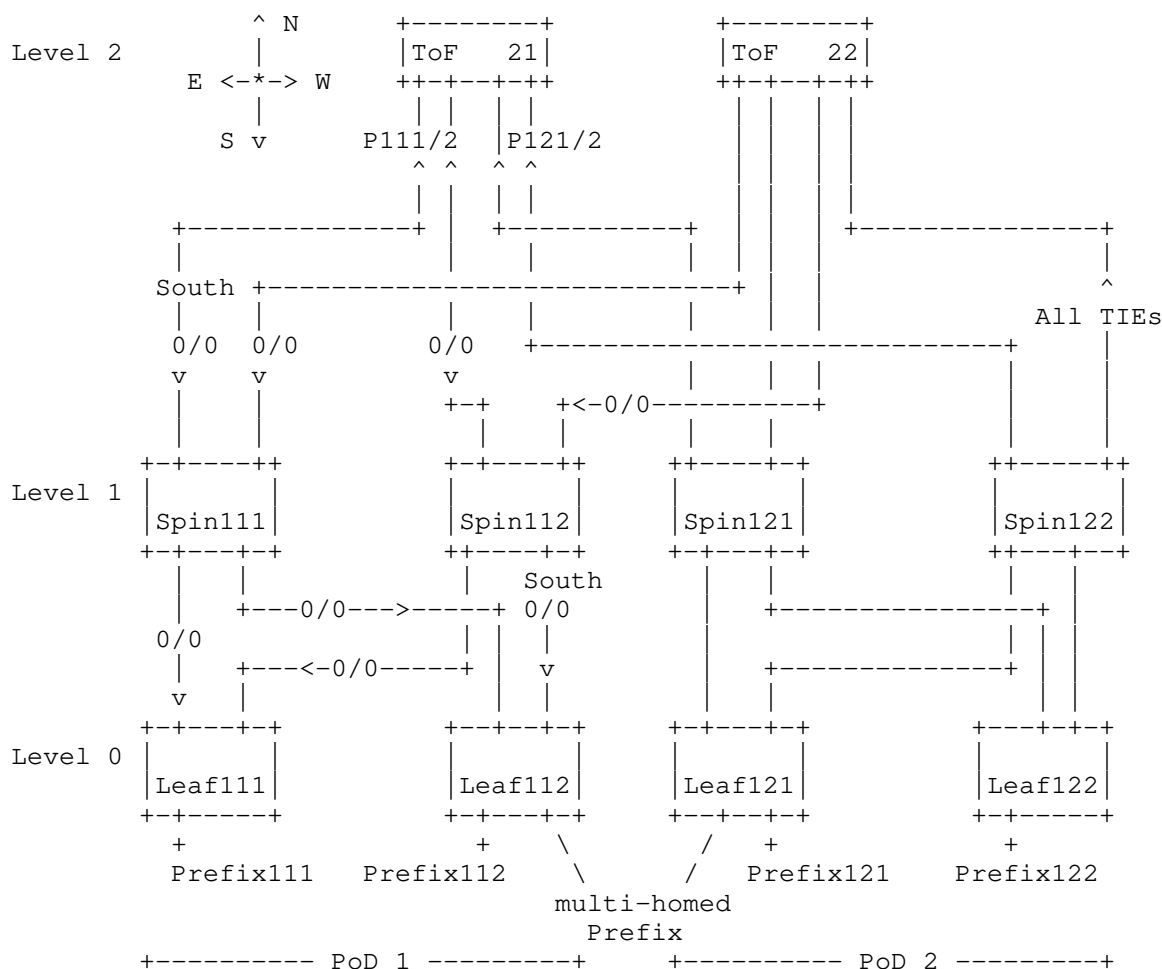


Figure 34: Normal Case Topology

This section describes RIFT deployment in example topology given in Figure 34 without any node or link failures. The scenario disregards flooding reduction for simplicity's sake and compresses 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

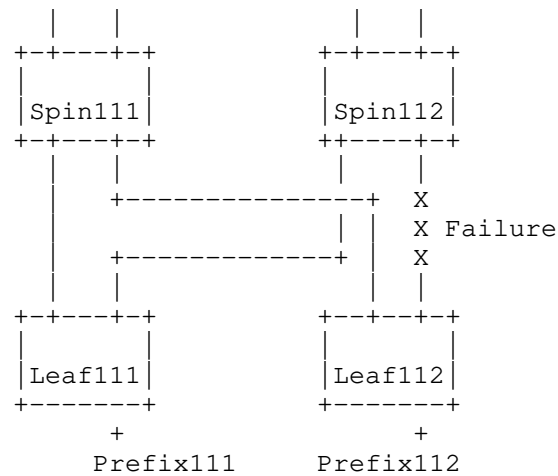


Figure 35: Single Leaf Link Failure

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 disaggregation is covered 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

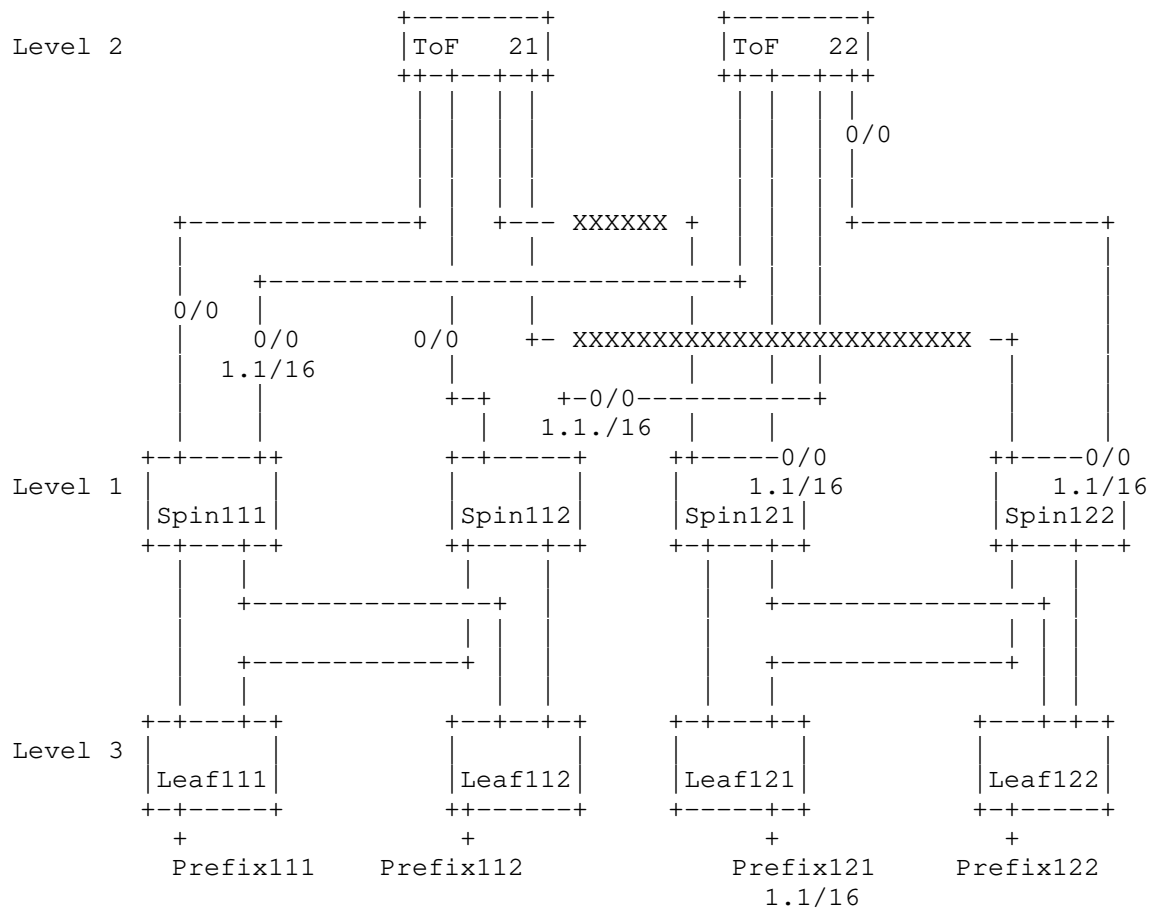


Figure 36: Fabric Partition

Figure 36 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 Spine 112 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 this 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

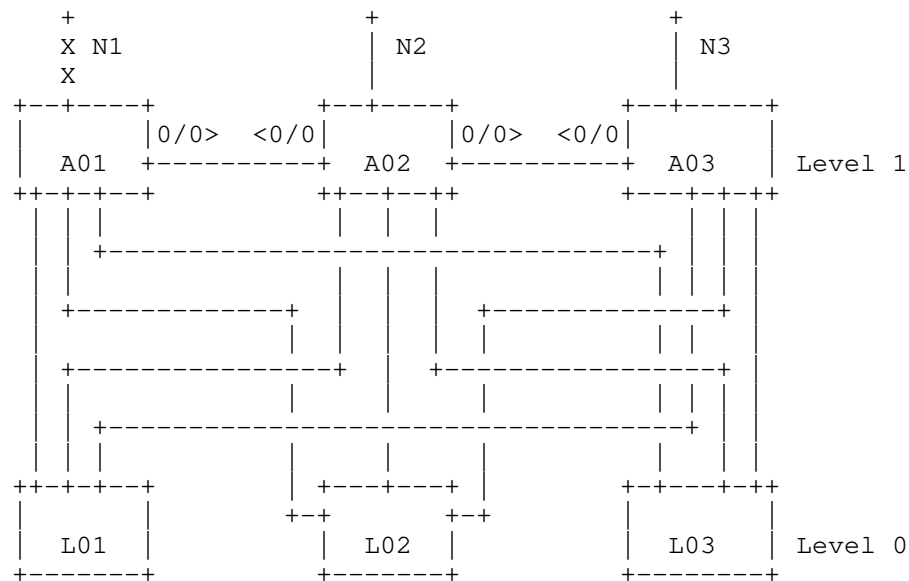


Figure 37: North Partitioned Router

Figure 37 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. Further Details on Implementation

### 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 disaggregation 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

Nodes that do not act as ToF are not required to discover fallen leaves by comparing reachable destinations with peers and therefore do not need to run the computation of disaggregated routes based on that discovery. On the other hand, non-ToF nodes need to respect disaggregated routes advertised from the north. In the case of negative disaggregation, spines nodes need to generate southbound disaggregated routes when all parents are lost for a fallen leaf.

## 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. Malformed Packets

The protocol protects packets extensively through optional signatures and nonces so if the possibility of maliciously injected malformed or replayed packets exist in a deployment, this conclusively protects against such attacks.

Even with security envelope, since RIFT relies on Thrift encoders and decoders generated automatically from IDL it is conceivable that errors in such encoders/decoders could be discovered and lead to delivery of corrupted packets or reception of packets that cannot be decoded. Misformatted packets lead normally to decoder returning an error condition to the caller and with that the packet is basically unparsable with no other choice but to discard it. Should the unlikely scenario occur of the decoder being forced to abort the protocol this is neither better nor worse than today's behavior of other protocols.

## 7.3. 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 ThreeWay 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.

## 7.4. 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.

## 7.5. 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.

#### 7.6. 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.

#### 7.7. 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.

#### 7.8. 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 the 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\_v5/common/AddressFamilyType"

Address family type.

## 8.2.1.1. Requested Entries

| Name                  | Value | Schema Version | Description |
|-----------------------|-------|----------------|-------------|
| Illegal               | 0     | 5.0            |             |
| AddressFamilyMinValue | 1     | 5.0            |             |
| IPv4                  | 2     | 5.0            |             |
| IPv6                  | 3     | 5.0            |             |
| AddressFamilyMaxValue | 4     | 5.0            |             |

Table 7

## 8.2.2. Registry RIFT\_v5/common/HierarchyIndications"

Flags indicating node configuration in case of ZTP.

## 8.2.2.1. Requested Entries

| Name                                 | Value | Schema Version | Description |
|--------------------------------------|-------|----------------|-------------|
| leaf_only                            | 0     | 5.0            |             |
| leaf_only_and_leaf_2_leaf_procedures | 1     | 5.0            |             |
| top_of_fabric                        | 2     | 5.0            |             |

Table 8

## 8.2.3. Registry RIFT\_v5/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     | 5.0            |             |
| AS_nsec | 2     | 5.0            |             |

Table 9

## 8.2.4. Registry RIFT\_v5/common/IPAddressType"

IP address type.

## 8.2.4.1. Requested Entries

| Name        | Value | Schema Version | Description     |
|-------------|-------|----------------|-----------------|
| ipv4address | 1     | 5.0            | Content is IPv4 |
| ipv6address | 2     | 5.0            | Content is IPv6 |

Table 10

## 8.2.5. Registry RIFT\_v5/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     | 5.0            |             |
| ipv6prefix | 2     | 5.0            |             |

Table 11



## 8.2.6. Registry RIFT\_v5/common/IPv4PrefixType"

IPv4 prefix type.

## 8.2.6.1. Requested Entries

| Name      | Value | Schema Version | Description |
|-----------|-------|----------------|-------------|
| address   | 1     | 5.0            |             |
| prefixlen | 2     | 5.0            |             |

Table 12

## 8.2.7. Registry RIFT\_v5/common/IPv6PrefixType"

IPv6 prefix type.

## 8.2.7.1. Requested Entries

| Name      | Value | Schema Version | Description |
|-----------|-------|----------------|-------------|
| address   | 1     | 5.0            |             |
| prefixlen | 2     | 5.0            |             |

Table 13

## 8.2.8. Registry RIFT\_v5/common/PrefixSequenceType"

Sequence of a prefix in case of move.

## 8.2.8.1. Requested Entries

| Name          | Value | Schema Version | Description                                      |
|---------------|-------|----------------|--------------------------------------------------|
| timestamp     | 1     | 5.0            |                                                  |
| transactionid | 2     | 5.0            | Transaction ID set by client in e.g. in 6LoWPAN. |

Table 14

## 8.2.9. Registry RIFT\_v5/common/RouteType"

RIFT route types. @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     | 5.0            |             |
| RouteTypeMinValue   | 1     | 5.0            |             |
| Discard             | 2     | 5.0            |             |
| LocalPrefix         | 3     | 5.0            |             |
| SouthPGPPrefix      | 4     | 5.0            |             |
| NorthPGPPrefix      | 5     | 5.0            |             |
| NorthPrefix         | 6     | 5.0            |             |
| NorthExternalPrefix | 7     | 5.0            |             |
| SouthPrefix         | 8     | 5.0            |             |
| SouthExternalPrefix | 9     | 5.0            |             |
| NegativeSouthPrefix | 10    | 5.0            |             |
| RouteTypeMaxValue   | 11    | 5.0            |             |

Table 15

## 8.2.10. Registry RIFT\_v5/common/TIETypeType"

Type of TIE.

## 8.2.10.1. Requested Entries

| Name                                        | Value | Schema<br>Version | Description |
|---------------------------------------------|-------|-------------------|-------------|
| Illegal                                     | 0     | 5.0               |             |
| TIETypeMinValue                             | 1     | 5.0               |             |
| NodeTIEType                                 | 2     | 5.0               |             |
| PrefixTIEType                               | 3     | 5.0               |             |
| PositiveDisaggregationPrefixTIEType         | 4     | 5.0               |             |
| NegativeDisaggregationPrefixTIEType         | 5     | 5.0               |             |
| PGFPrefixTIEType                            | 6     | 5.0               |             |
| KeyValueTIEType                             | 7     | 5.0               |             |
| ExternalPrefixTIEType                       | 8     | 5.0               |             |
| PositiveExternalDisaggregationPrefixTIEType | 9     | 5.0               |             |
| TIETypeMaxValue                             | 10    | 5.0               |             |

Table 16

## 8.2.11. Registry RIFT\_v5/common/TieDirectionType"

Direction of TIEs.

## 8.2.11.1. Requested Entries

| Name              | Value | Schema Version | Description |
|-------------------|-------|----------------|-------------|
| Illegal           | 0     | 5.0            |             |
| South             | 1     | 5.0            |             |
| North             | 2     | 5.0            |             |
| DirectionMaxValue | 3     | 5.0            |             |

Table 17

## 8.2.12. Registry RIFT\_v5/encoding/Community"

Prefix community.

## 8.2.12.1. Requested Entries

| Name   | Value | Schema Version | Description       |
|--------|-------|----------------|-------------------|
| top    | 1     | 5.0            | Higher order bits |
| bottom | 2     | 5.0            | Lower order bits  |

Table 18

## 8.2.13. Registry RIFT\_v5/encoding/KeyValueTIEElement"

Generic key value pairs.

## 8.2.13.1. Requested Entries

| Name      | Value | Schema Version | Description |
|-----------|-------|----------------|-------------|
| keyvalues | 1     | 5.0            |             |

Table 19

## 8.2.14. Registry RIFT\_v5/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 Version | Description                                                        |
|-------------------|-------|----------------|--------------------------------------------------------------------|
| name              | 1     | 5.0            | Node or adjacency name.                                            |
| local_id          | 2     | 5.0            | Local link ID.                                                     |
| flood_port        | 3     | 5.0            | UDP port to which we can receive flooded TIEs.                     |
| link_mtu_size     | 4     | 5.0            | Layer 3 MTU, used to discover mismatch.                            |
| link_bandwidth    | 5     | 5.0            | Local link bandwidth on the interface.                             |
| neighbor          | 6     | 5.0            | Reflects the neighbor once received to provide 3-way connectivity. |
| pod               | 7     | 5.0            | Node's PoD.                                                        |
| node_capabilities | 10    | 5.0            | Node capabilities supported.                                       |
| link_capabilities | 11    | 5.0            | Capabilities of this link.                                         |
| holdtime          | 12    | 5.0            | Required holdtime of the                                           |

|                             |    |     |                                                                                                                                                              |
|-----------------------------|----|-----|--------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                             |    |     | adjacency, i.e. for how long a period should adjacency be kept up without valid LIE reception.                                                               |
| label                       | 13 | 5.0 | Optional, unsolicited, downstream assigned locally significant label value for the adjacency.                                                                |
| not_a_ztp_offer             | 21 | 5.0 | 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 | 5.0 | Indicates to northbound neighbor that it should be reflooding TIEs received from this node to achieve flood reduction and balancing for northbound flooding. |
| you_are_sending_too_quickly | 23 | 5.0 | Indicates to neighbor to flood node TIEs only and slow down all other TIEs. Ignored when received from southbound neighbor.                                  |

|               |    |     |                                                                          |
|---------------|----|-----|--------------------------------------------------------------------------|
| instance_name | 24 | 5.0 | Instance name in case multiple RIFT instances running on same interface. |
|---------------|----|-----|--------------------------------------------------------------------------|

Table 20

## 8.2.15. Registry RIFT\_v5/encoding/LinkCapabilities"

Link capabilities.

## 8.2.15.1. Requested Entries

| Name                    | Value | Schema Version | Description                                                   |
|-------------------------|-------|----------------|---------------------------------------------------------------|
| bfd                     | 1     | 5.0            | Indicates that the link is supporting BFD.                    |
| ipv4_forwarding_capable | 2     | 5.0            | Indicates whether the interface will support IPv4 forwarding. |

Table 21

## 8.2.16. Registry RIFT\_v5/encoding/LinkIDPair"

LinkID pair describes one of parallel links between two nodes.

## 8.2.16.1. Requested Entries

| Name      | Value | Schema Version | Description                                |
|-----------|-------|----------------|--------------------------------------------|
| local_id  | 1     | 5.0            | Node-wide unique value for the local link. |
| remote_id | 2     | 5.0            | Received remote link ID for this           |

|                            |    |     |                                                                                                                                                         |
|----------------------------|----|-----|---------------------------------------------------------------------------------------------------------------------------------------------------------|
|                            |    |     | link.                                                                                                                                                   |
| platform_interface_index   | 10 | 5.0 | Describes the local interface index of the link.                                                                                                        |
| platform_interface_name    | 11 | 5.0 | Describes the local interface name.                                                                                                                     |
| trusted_outer_security_key | 12 | 5.0 | Indicates 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 | 5.0 | Indicates whether the link is protected by established BFD session.                                                                                     |
| address_families           | 14 | 5.0 | Optional indication which address families are up on the interface                                                                                      |

Table 22

## 8.2.17. Registry RIFT\_v5/encoding/Neighbor"

Neighbor structure.



## 8.2.17.1. Requested Entries

| Name       | Value | Schema Version | Description                    |
|------------|-------|----------------|--------------------------------|
| originator | 1     | 5.0            | System ID of the originator.   |
| remote_id  | 2     | 5.0            | ID of remote side of the link. |

Table 23

## 8.2.18. Registry RIFT\_v5/encoding/NodeCapabilities"

Capabilities the node supports.

## 8.2.18.1. Requested Entries

| Name                   | Value | Schema Version | Description                                                                                                   |
|------------------------|-------|----------------|---------------------------------------------------------------------------------------------------------------|
| protocol_minor_version | 1     | 5.0            | Must advertise supported minor version dialect that way.                                                      |
| flood_reduction        | 2     | 5.0            | indicates that node supports flood reduction.                                                                 |
| hierarchy_indications  | 3     | 5.0            | indicates place in hierarchy, i.e. top-of-fabric or leaf only (in ZTP) or support for leaf-2-leaf procedures. |

Table 24

## 8.2.19. Registry RIFT\_v5/encoding/NodeFlags"

Indication flags of the node.

## 8.2.19.1. Requested Entries

| Name     | Value | Schema Version | Description                                                            |
|----------|-------|----------------|------------------------------------------------------------------------|
| overload | 1     | 5.0            | Indicates that node is in overload, do not transit traffic through it. |

Table 25

## 8.2.20. Registry RIFT\_v5/encoding/NodeNeighborsTIEElement"

neighbor of a node

## 8.2.20.1. Requested Entries

| Name      | Value | Schema Version | Description                                                                              |
|-----------|-------|----------------|------------------------------------------------------------------------------------------|
| level     | 1     | 5.0            | level of neighbor                                                                        |
| cost      | 3     | 5.0            | Cost to neighbor. Ignore anything larger than 'infinite_distance' and 'invalid_distance' |
| link_ids  | 4     | 5.0            | can carry description of multiple parallel links in a TIE                                |
| bandwidth | 5     | 5.0            | total bandwidth to neighbor as sum of all parallel links                                 |

Table 26

## 8.2.21. Registry RIFT\_v5/encoding/NodeTIEElement"

Description of a node.

## 8.2.21.1. Requested Entries

| Name            | Value | Schema Version | Description                                                                |
|-----------------|-------|----------------|----------------------------------------------------------------------------|
| level           | 1     | 5.0            | Level of the node.                                                         |
| neighbors       | 2     | 5.0            | Node's neighbors. Multiple node TIEs can carry disjoint sets of neighbors. |
| capabilities    | 3     | 5.0            | Capabilities of the node.                                                  |
| flags           | 4     | 5.0            | Flags of the node.                                                         |
| name            | 5     | 5.0            | Optional node name for easier operations.                                  |
| pod             | 6     | 5.0            | PoD to which the node belongs.                                             |
| startup_time    | 7     | 5.0            | optional startup time of the node                                          |
| miscabled_links | 10    | 5.0            | If any local links are miscabled, this indication is flooded.              |

Table 27

## 8.2.22. Registry RIFT\_v5/encoding/PacketContent"

Content of a RIFT packet.

## 8.2.22.1. Requested Entries

| Name | Value | Schema Version | Description |
|------|-------|----------------|-------------|
| lie  | 1     | 5.0            |             |
| tide | 2     | 5.0            |             |
| tire | 3     | 5.0            |             |
| tie  | 4     | 5.0            |             |

+-----+-----+-----+-----+

Table 28

## 8.2.23. Registry RIFT\_v5/encoding/PacketHeader"

Common RIFT packet header.

## 8.2.23.1. Requested Entries

| Name          | Value | Schema<br>Version | Description                                                                                                                                                                |
|---------------|-------|-------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| major_version | 1     | 5.0               | Major version of protocol.                                                                                                                                                 |
| minor_version | 2     | 5.0               | Minor version of protocol.                                                                                                                                                 |
| sender        | 3     | 5.0               | Node sending the packet, in<br>case of LIE/TIRE/TIDE also<br>the originator of it.                                                                                         |
| level         | 4     | 5.0               | Level of the node sending the<br>packet, required on<br>everything except LIEs. Lack<br>of presence on LIEs indicates<br>UNDEFINED_LEVEL and is used<br>in ZTP procedures. |

Table 29

## 8.2.24. Registry RIFT\_v5/encoding/PrefixAttributes"

Attributes of a prefix.

## 8.2.24.1. Requested Entries

| Name   | Value | Schema<br>Version | Description                                                              |
|--------|-------|-------------------|--------------------------------------------------------------------------|
| metric | 2     | 5.0               | Distance of the<br>prefix.                                               |
| tags   | 3     | 5.0               | Generic unordered set<br>of route tags, can be<br>redistributed to other |

|                   |    |     |                                                             |
|-------------------|----|-----|-------------------------------------------------------------|
|                   |    |     | protocols or use within the context of real time analytics. |
| monotonic_clock   | 4  | 5.0 | Monotonic clock for mobile addresses.                       |
| loopback          | 6  | 5.0 | Indicates if the prefix is a node loopback.                 |
| directly_attached | 7  | 5.0 | Indicates that the prefix is directly attached.             |
| from_link         | 10 | 5.0 | link to which the address belongs to.                       |

Table 30

## 8.2.25. Registry RIFT\_v5/encoding/PrefixTIEElement"

TIE carrying prefixes

## 8.2.25.1. Requested Entries

| Name     | Value | Schema Version | Description                              |
|----------|-------|----------------|------------------------------------------|
| prefixes | 1     | 5.0            | Prefixes with the associated attributes. |

Table 31

## 8.2.26. Registry RIFT\_v5/encoding/ProtocolPacket"

RIFT packet structure.

## 8.2.26.1. Requested Entries

| Name    | Value | Schema Version | Description |
|---------|-------|----------------|-------------|
| header  | 1     | 5.0            |             |
| content | 2     | 5.0            |             |

+-----+-----+-----+-----+

Table 32

## 8.2.27. Registry RIFT\_v5/encoding/TIDEPacket"

TIDE with \*sorted\* TIE headers.

## 8.2.27.1. Requested Entries

| Name        | Value | Schema Version | Description                          |
|-------------|-------|----------------|--------------------------------------|
| start_range | 1     | 5.0            | First TIE header in the tide packet. |
| end_range   | 2     | 5.0            | Last TIE header in the tide packet.  |
| headers     | 3     | 5.0            | _Sorted_ list of headers.            |

Table 33

## 8.2.28. Registry RIFT\_v5/encoding/TIEElement"

Single element in a TIE.

## 8.2.28.1. Requested Entries

| Name                                                  | Value | Schema Version | Description                   |
|-------------------------------------------------------|-------|----------------|-------------------------------|
| node case of enum NodeTIEType.                        | 1     | 5.0            | Used in common.TIETypeType.   |
| prefixes case of enum PrefixTIEType.                  | 2     | 5.0            | Used in common.TIETypeType.Pr |
| positive_disaggregation_prefixes (always southbound). | 3     | 5.0            | Positive pref                 |
| negative_disaggregation_prefixes                      | 5     | 5.0            | Transitive, negat             |



|                                           |   |     |                     |
|-------------------------------------------|---|-----|---------------------|
| positive_external_disaggregation_prefixes | 7 | 5.0 | Positive external d |
| isaggregated                              |   |     | prefixes (always    |
| southbound).                              |   |     |                     |
| keyvalues                                 | 9 | 5.0 | Key-Value sto       |
| re elements.                              |   |     |                     |

Table 34

## 8.2.29. Registry RIFT\_v5/encoding/TIEHeader"

Header of a TIE.

## 8.2.29.1. Requested Entries

| Name                 | Value | Schema Version | Description                                    |
|----------------------|-------|----------------|------------------------------------------------|
| tieid                | 2     | 5.0            | ID of the tie.                                 |
| seq_nr               | 3     | 5.0            | Sequence number of the tie.                    |
| origination_time     | 10    | 5.0            | Absolute timestamp when the TIE was generated. |
| origination_lifetime | 12    | 5.0            | Original lifetime when the TIE was generated.  |

Table 35

## 8.2.30. Registry RIFT\_v5/encoding/TIEHeaderWithLifeTime"

Header of a TIE as described in TIRE/TIDE.





## 8.2.30.1. Requested Entries

| Name               | Value | Schema Version | Description         |
|--------------------|-------|----------------|---------------------|
| header             | 1     | 5.0            |                     |
| remaining_lifetime | 2     | 5.0            | Remaining lifetime. |

Table 36

## 8.2.31. Registry RIFT\_v5/encoding/TIEID"

Unique ID of a TIE.

## 8.2.31.1. Requested Entries

| Name       | Value | Schema Version | Description                     |
|------------|-------|----------------|---------------------------------|
| direction  | 1     | 5.0            | direction of TIE                |
| originator | 2     | 5.0            | indicates originator of the TIE |
| tietype    | 3     | 5.0            | type of the tie                 |
| tie_nr     | 4     | 5.0            | number of the tie               |

Table 37

## 8.2.32. Registry RIFT\_v5/encoding/TIEPacket"

TIE packet

## 8.2.32.1. Requested Entries

| Name    | Value | Schema Version | Description |
|---------|-------|----------------|-------------|
| header  | 1     | 5.0            |             |
| element | 2     | 5.0            |             |

Table 38

## 8.2.33. Registry RIFT\_v5/encoding/TIREPacket"

TIRE packet

## 8.2.33.1. Requested Entries

| Name    | Value | Schema Version | Description |
|---------|-------|----------------|-------------|
| headers | 1     | 5.0            |             |

Table 39

## 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. Contributors

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. |       | Alankar Sharma |       | Pascal Thubert   |
| -----                | ----- | -----          | ----- | -----            |
| Juniper              |       | Comcast        |       | Cisco            |
| -----                | ----- | -----          | ----- | -----            |
| Bruno Rijsman        |       | Jordan Head    |       | Dmitry Afanasiev |
| -----                | ----- | -----          | ----- | -----            |
| Individual           |       | Juniper        |       | Yandex           |
| -----                | ----- | -----          | ----- | -----            |
| Don Fedyk            |       | Alia Atlas     |       | John Drake       |
| -----                | ----- | -----          | ----- | -----            |
| Individual           |       | Individual     |       | Juniper          |
| -----                | ----- | -----          | ----- | -----            |
| Ilya Vershkov        |       |                |       |                  |
| -----                | ----- | -----          | ----- | -----            |
| Mellanox             |       |                |       |                  |
| -----                | ----- | -----          | ----- | -----            |

Table 40: RIFT Authors

## 11. References

### 11.1. 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/develo/regist/tut/eui.pdf>>.
- [RFC1982] Elz, R. and R. Bush, "Serial Number Arithmetic", RFC 1982, DOI 10.17487/RFC1982, August 1996, <<https://www.rfc-editor.org/info/rfc1982>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2365] Meyer, D., "Administratively Scoped IP Multicast", BCP 23, RFC 2365, DOI 10.17487/RFC2365, July 1998, <<https://www.rfc-editor.org/info/rfc2365>>.
- [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-IS)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.
- [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>>.
- [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>>.
- [VFR] Erlebach et al., T., "Cuts and Disjoint Paths in the Valley-Free Path Model of Internet BGP Routing", Springer Berlin Heidelberg Combinatorial and Algorithmic Aspects of Networking, 2005.

## 11.2. Informative References

- [APPLICABILITY] Wei, Y., Zhang, Z., Afanasiev, D., Thubert, P., Verhaeg, T., and J. Kowalczyk, "RIFT Applicability", Work in Progress, Internet-Draft, draft-ietf-rift-applicability-05, 26 April 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-rift-applicability-05>>.

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

- [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 the 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.TITypeType'` (use of enumeration types is generally discouraged) *\*or\**
11. add new TIE type to `'TITypeType'` with flooding scope different from prefix TIE flooding scope

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

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. Backwards-Compatible Extension of Schema

The set of rules in Appendix B 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'. Contrary to that, 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.

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. However, it is important to understand that a higher minor version of a schema does *\*not\** guarantee that capabilities introduced in lower minors of the same major are supported. The `'node_capabilities'` field is used to indicate which capabilities are supported.

Specifically, the schema may add elements to `'NodeCapabilities'` field future capabilities to indicate whether it will support interpretation of 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 will be however unavoidable. `'NodeCapabilities'`

shown in LIE MUST match the capabilities shown in the Node TIEs, otherwise the behavior is unspecified. A node detecting the mismatch SHOULD generate a notification.

To support new TIE types without increasing the major version enumeration 'TIEElement' can be extended with new optional elements for new 'common.TIETypeType' values as long the scope of the new TIE matches the prefix TIE scope. In case it is necessary to understand whether all nodes can parse the new TIE type a node capability MUST be added in 'NodeCapabilities' to prevent a non-homogenous network.

## B.2. common.thrift

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

namespace py common

/** @note MUST be interpreted in implementation as unsigned 64 bits.
 */
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
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 hence 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 i32      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
/** @note MUST be interpreted in implementation as unsigned,
    especially since we have the /128 IPv6 case. */
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
/** used when node is configured as top of fabric in ZTP.*/
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

/** invalid key for key value */
const KeyIDType invalid_key_value_key = 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 LIE TX interval time */
const TimeIntervalInSecType default_lie_tx_interval = 1
/** default LIE FSM holddown time */
const TimeIntervalInSecType default_lie_holdtime = 3
/** multiplier for default_lie_holdtime to hold down multiple neighbors */
const i8 multiple_neighbors_lie_holdtime_multiplifier = 4
/** 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
/** optional 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). */
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. */
const i16 maximum_valid_nonce_delta = 5;
const TimeIntervalInSecType nonce_regeneration_interval = 300;

/** 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;
} (python.immutable = "")

/** IPv6 prefix type. */
struct IPv6PrefixType {
    1: required IPv6Address    address;
    2: required PrefixLenType  prefixlen;
} (python.immutable = "")

/** IP address type. */
union IPAddressType {
    /** Content is IPv4 */
    1: optional IPv4Address    ipv4address;
    /** Content is IPv6 */
    2: optional IPv6Address    ipv6address;
} (python.immutable = "")

/** 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;
} (python.immutable = "")

/** 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.
*/
enum TIETypeType {
    Illegal                                = 0,
    TIETypeMinValue                        = 1,
    /** first legal value */
    NodeTIEType                            = 2,

```

```
    PrefixTIEType                = 3,
    PositiveDisaggregationPrefixTIEType = 4,
    NegativeDisaggregationPrefixTIEType = 5,
    PGPPrefixTIEType              = 6,
    KeyValueTIEType               = 7,
    ExternalPrefixTIEType         = 8,
    PositiveExternalDisaggregationPrefixTIEType = 9,
    TIETypeMaxValue               = 10,
}

/** RIFT route types.
    @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.3. encoding.thrift



```
/**
    Thrift file for packet encodings for RIFT
*/

include "common.thrift"

namespace py encoding

/** Represents protocol encoding schema major version */
const common.VersionType protocol_major_version = 5
/** Represents protocol encoding schema minor version */
const common.MinorVersionType protocol_minor_version = 0

/** 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;
} (python.immutable = "")

/** 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;
} (python.immutable = "")
```

```
/** Capabilities the node supports. */
struct NodeCapabilities {
    /** Must advertise supported minor version dialect that way. */
    1: required common.MinorVersionType      protocol_minor_version =
        protocol_minor_version;
    /** indicates that node supports flood reduction. */
    2: optional bool                        flood_reduction =
        common.flood_reduction_default;
    /** indicates place in hierarchy, i.e. top-of-fabric or
        leaf only (in ZTP) or support for leaf-2-leaf
        procedures. */
    3: optional common.HierarchyIndications  hierarchy_indications;
} (python.immutable = "")

/** Link capabilities. */
struct LinkCapabilities {
    /** Indicates that the link is supporting BFD. */
    1: optional bool                        bfd =
        common.bfd_default;
    /** Indicates whether the interface will support IPv4 forwarding. */
    2: optional bool                        ipv4_forwarding_capable =
        true;
} (python.immutable = "")

/** 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 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 supported. */
10: required NodeCapabilities          node_capabilities;
    /** Capabilities of this link. */
11: optional LinkCapabilities          link_capabilities;
    /** Required holdtime of the adjacency, i.e. for how
        long a period should adjacency be kept up without valid LIE reception. */
12: required common.TimeIntervalInSecType
        holdtime = common.default_lie_holdtime;
    /** Optional, 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 reflooding TIEs received from this node to achieve flood
        reduction and balancing for northbound flooding. */
22: optional bool                      you_are_flood_repeater =
        common.default_you_are_flood_repeater;
    /** Indicates to neighbor to flood node TIEs only and slow down
        all other TIEs. Ignored when received from southbound neighbor. */
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;
    /** Indicates 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;
    /** Indicates 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>
        (python.immutable = "") address_families;
} (python.immutable = "")

/** Unique ID of a TIE. */
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;
} (python.immutable = "")

/** Header of a TIE. */
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. */
    10: optional common.IEEE802_1ASTimeStampType origination_time;
    /** Original lifetime when the TIE was generated. */
    12: optional common.LifeTimeInSecType origination_lifetime;
}

/** Header of a TIE as described in TIRE/TIDE.
 */
struct TIEHeaderWithLifeTime {
    1: required TIEHeader header;
    /** Remaining lifetime. */
    2: required common.LifeTimeInSecType remaining_lifetime;
}

/** TIDE with *sorted* TIE headers. */
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>
        (python.immutable = "") headers;
}

```

```

}

/** TIRE packet */
struct TIREPacket {
    1: required set<TIEHeaderWithLifeTime>
        (python.immutable = "") headers;
}

/** neighbor of a node */
struct NodeNeighborsTIEElement {
    /** level of neighbor */
    1: required common.LevelType level;
    /** Cost to neighbor. Ignore anything larger than 'infinite_distance' and 'i
nvalid_distance' */
    3: optional common.MetricType cost
        = common.default_distance;
    /** can carry description of multiple parallel links in a TIE */
    4: optional set<LinkIDPair>
        (python.immutable = "") link_ids;
    /** total bandwidth to neighbor as sum of all parallel links */
    5: optional common.BandwidthInMegaBitsType
        bandwidth = common.default_bandwidth;
} (python.immutable = "")

/** 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;
} (python.immutable = "")

/** Description of a node. */
struct NodeTIEElement {
    /** Level of the node. */
    1: required common.LevelType level;
    /** Node's neighbors. Multiple node TIEs can carry disjoint sets of neighbors
    . */
    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, this indication is flooded. */
    10: optional set<common.LinkIDType>
        (python.immutable = "")                miscabled_links;

} (python.immutable = "")

/** 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>
        (python.immutable = "")                tags;
    /** Monotonic clock for mobile addresses. */
    4: optional common.PrefixSequenceType        monotonic_clock;
    /** Indicates if the prefix is a node loopback. */
    6: optional bool                            loopback = false;
    /** Indicates that the prefix is directly attached. */
    7: optional bool                            directly_attached = true;
    /** link to which the address belongs to. */
    10: optional common.LinkIDType              from_link;
} (python.immutable = "")

/** TIE carrying prefixes */
struct PrefixTIEElement {
    /** Prefixes with the associated attributes. */
    1: required map<common.IPPrefixType, PrefixAttributes> prefixes;
} (python.immutable = "")

/** Generic key value pairs. */
struct KeyValueTIEElement {
    1: required map<common.KeyIDType, binary>    keyvalues;
} (python.immutable = "")

/** Single element in a TIE. */
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). */
    3: optional PrefixTIEElement              positive_disaggregation_prefixes;
    /** Transitive, negative prefixes (always southbound) */
    5: optional PrefixTIEElement              negative_disaggregation_prefixes;
    /** Externally reimported prefixes. */

```

```
6: optional PrefixTIEElement      external_prefixes;
/** Positive external disaggregated prefixes (always southbound). */
7: optional PrefixTIEElement
    positive_external_disaggregation_prefixes;
/** Key-Value store elements. */
9: optional KeyValueTIEElement keyvalues;
} (python.immutable = "")

/** 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=TIETTypeMinValue, tie_nr=0, direction=South)                   |
| MAX_TIEID signifies end of TIDEs   | Constant                    | TIE Key with maximal values: TIEID(originator=MAX_UINT64, tietype=TIETTypeMaxValue, tie_nr=MAX_UINT64, direction=North) |

Table 41: all\_constants



Authors' Addresses

Tony Przygienda (editor)  
Juniper  
1137 Innovation Way  
Sunnyvale, CA  
United States of America

Email: prz@juniper.net

Alankar Sharma  
Comcast  
1800 Bishops Gate Blvd  
Mount Laurel, NJ 08054  
United States of America

Email: Alankar\_Sharma@comcast.com

Pascal Thubert  
Cisco Systems, Inc  
Building D  
45 Allee des Ormes - BP1200  
06254 MOUGINS - Sophia Antipolis  
France

Phone: +33 497 23 26 34  
Email: pthubert@cisco.com

Bruno Rijsman  
Individual

Email: brunorijsman@gmail.com

Dmitry Afanasiev  
Yandex

Email: fl0w@yandex-team.ru

Internet Engineering Task Force  
Internet-Draft  
Intended status: Standards Track  
Expires: 8 January 2022

Z. Zhang  
Y. Wei  
B. Xu  
ZTE Corporation  
7 July 2021

Supporting leaves without northbound neighbors connecting to a fat-tree  
network using RIFT  
draft-zwx-rift-leaf-ring-00

#### Abstract

This document discusses the usage and solution for leaf nodes without northbound neighbors connecting to a fat-tree network by leaf nodes having direct northbound neighbors in RIFT.

#### 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 8 January 2022.

#### Copyright Notice

Copyright (c) 2021 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 |
| 1.1. Requirements Language . . . . .    | 2 |
| 2. Problem Statement . . . . .          | 2 |
| 3. Solution . . . . .                   | 6 |
| 3.1. Capability advertisement . . . . . | 6 |
| 3.2. Prefix transferring . . . . .      | 6 |
| 3.3. Example . . . . .                  | 7 |
| 4. IANA Considerations . . . . .        | 7 |
| 5. Security Considerations . . . . .    | 7 |
| 6. References . . . . .                 | 7 |
| 6.1. Normative References . . . . .     | 7 |
| 6.2. Informative References . . . . .   | 7 |
| Authors' Addresses . . . . .            | 8 |

## 1. Introduction

[I-D.ietf-rift-rift] specifies a dynamic routing protocol for Clos and fat-tree network topology. It suits most of the deployments. In some situations, the leaves are connected by ring or chain topology, and some of them may have no northbound link, so these nodes may not be able to connecting the network. This document discusses the usage and proposes a solution.

## 1.1. 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 2119 [RFC2119].

## 2. Problem Statement

[I-D.ietf-rift-rift] specifies a dynamic routing protocol for Clos and fat-tree network topology. The leaf part of a traditional Clos and fat-tree network topology is shown as:

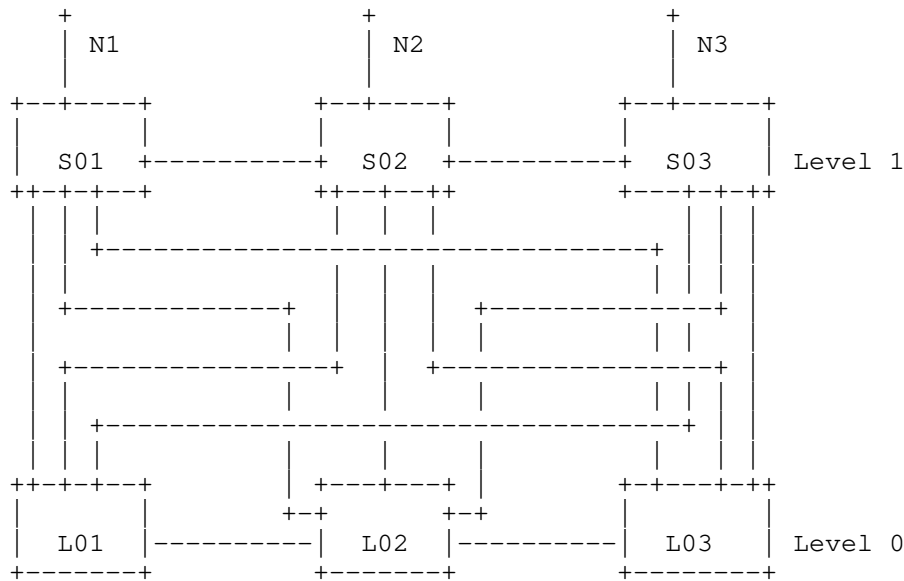


Figure 1

In most of the cases, each leaf node has north connections with at least one spine, and there may be east-west links between leaf nodes. In case a leaf node lost all the north connections, it can still access the network through the east-west link between leaves. For example in Figure 1, there is an east-west link between L01 and L02, and there is an east-west link between L02 and L03. When the northbound connections for the leaves are all work well, L01, L02 and L03 may generate a Prefix South TIE with default route and advertise it through the east-west links according to the definition in section 4.2.3.4 in [I-D.ietf-rist-rist]. In case L01 lost all the northbound links with S01, S02 and S03, according to Northbound SPF algorithm defined in section 4.2.4.1 in [I-D.ietf-rist-rist], L01 can compute the next hop L02 for the default route. On the other hand, the prefix of L01 can be flood northbound by L02. Then L01 can still access the network through the east-west link between L01 and L02.

But in some deployments, the leaves may connect with each other by ring topology (sometimes, a chain topology), and not all of them have northbound connection with spine nodes.

For example, in the IP Radio Access Network (IP RAN, mobile backhaul network), the 4G eNB or the 5G gNB connect to an IP access network of a ring topology. The access network attaches to an aggregation network through two aggregation nodes. In 5G era, the aggregation network and the metro network is evolving to Spine-and-Leaf

architecture to take the advantage of Spine-and-Leaf. Figure 2 depicts an diagram of an IPRAN network with Spine-and-Leaf architecture. If the aggregation network runs RIFT, using the proposal of this draft, the access network ring does not need to deploy other IGP protocol to enable the routing.

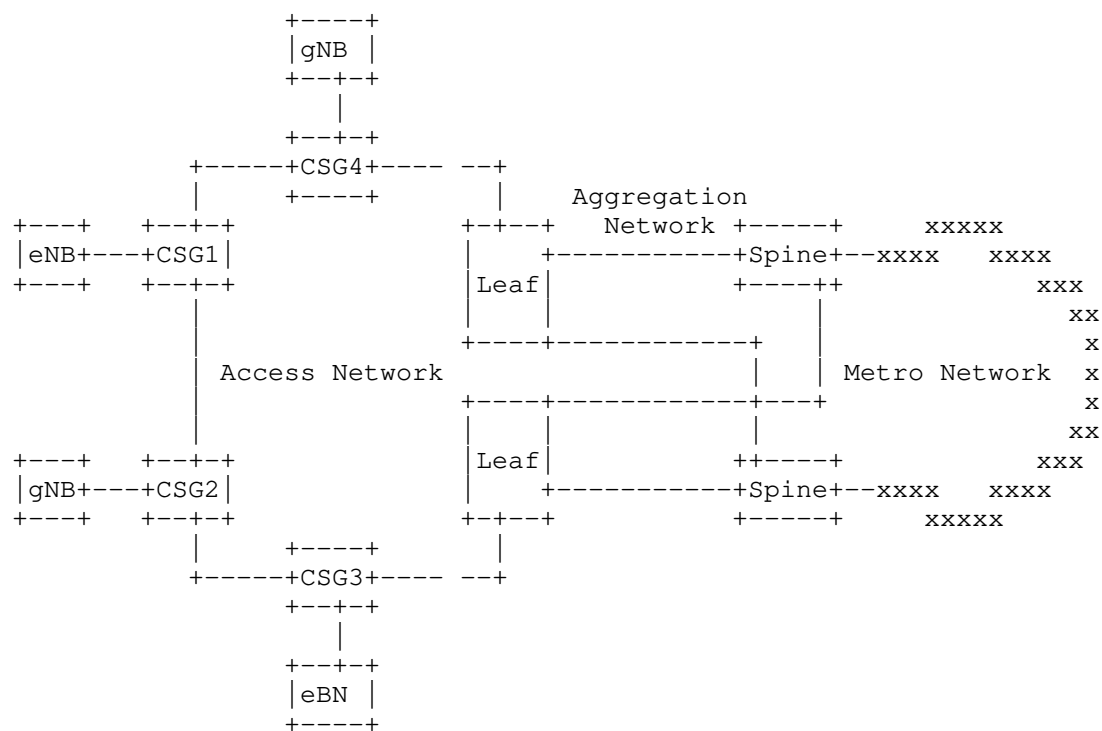


Figure 2: IPRAN network with Spine-and-Leaf architecture

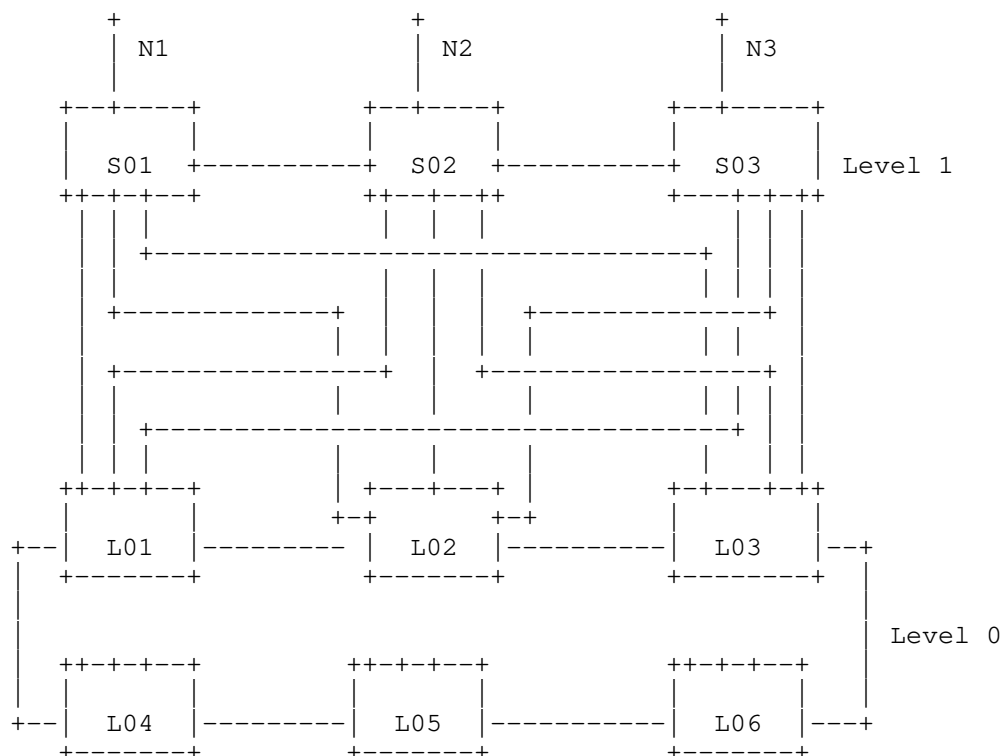


Figure 3

Figure 3 is an abstract of Figure 2, several leaves connect to the fat-tree network by ring topology, each leaf has two east-west connections with other leaves, but only L01, L02 and L03 have northbound connection with spine nodes. L01, L02 and L03 advertise a Prefix South TIE with default route through the east-west links. L04 and L06 can access the network through L01 and L03. But L04 and L06 do not generate or flood the Prefix South TIE with default route because they have no northbound link. So L05 cannot receive the Prefix South TIE with default route through the link between L04 and L05, or the link between L05 and L06. On the other hand, the prefix of L05 cannot be flooded east-west by L04 and L05. So L05 cannot send flow to other leaf, and other leaf cannot send flow to L05 also. L05 cannot access the network.

This document discuss the extension that can be used to solve the problem when some leaves without northbound neighbors connecting to a fat-tree network.

### 3. Solution

#### 3.1. Capability advertisement

A new link capability which is named leaf-transport is set in LIE and node TIE. The new capability indicates that the leaf node can transfer the Prefix TIE received from the east-west link to the other leaf node.

The LIE FSM will not be affected if only one leaf supports the capability and the neighbor does not support.

The capability can be used for diagnosis in case there is something wrong when computing route.

#### 3.2. Prefix transferring

When a leaf node which has no northbound connection receives Prefix TIE from a neighbor by an east-west link, it can transfer the Prefix TIE to the other neighbor by the other east-west link, with increased metric. The increased metric should be the sum of the received metric and the metric of the east-west link.

The Prefix TIE can be the default route and the prefix of the leaf node, other prefixes can also be transferred. The network administrator can control the prefix transferring by policy.

##### Prefix South TIE transferring

The leaf node without northbound neighbors which supports the leaf transfer capability, MUST transfer the Prefix South TIE received from an east-west neighbor to the other east-west neighbor which also has no northbound connection.

##### Prefix North TIE transferring

The leaf node without northbound neighbors which supports the leaf transfer capability, MUST transfer the Prefix North TIE received from an east-west neighbor which has no northbound connection to the other east-west neighbor.

The leaf node without northbound neighbors which supports the leaf transfer capability, MAY transfer the Prefix TIE from an east-west neighbor to the other east-west neighbor which has northbound connection. According to Northbound SPF algorithm defined in section 4.2.4.1 in [I-D.ietf-rift-rift], the transferring does not affect the routing calculating result of the neighbors which has northbound connection.

### 3.3. Example

In figure 2, either L04 or L06, or both of them advertise the leaf-transport capability in LIE to L05 when they are forming an adjacency. And the leaf-transport capability is also set in node TIE when L04 or L06 advertise the TIE.

L04/L06 receives Prefix South TIE with default route from L01/L03, L04/L06 transfers the route through Prefix South TIE with increased metric to L05.

L04/L06 receives Prefix North TIE of L05 and transfers the route through Prefix North TIE with increased metric to L01/L03.

L05 receives the Prefix South TIE from L04/L06, after N-SPF calculation, L05 calculates the default route with next hop L04/L06. L01/L03 receives prefix of L05 from L04/L06, and L01/L03 floods the Prefix North TIE northbound. Then L05 can send flow to other leaf through L04/L06. The flow sent by other leaf with destination set to the prefix of L05 can also be routed to L05. Then L05 can access the network.

## 4. IANA Considerations

A new type for 'leaf-transfer' is requested in link capability.

## 5. Security Considerations

When the node without northbound neighbors supports the function defined in this document, there may be unnecessary Prefix TIE advertisement, and unnecessary prefix may be leaked.

## 6. References

### 6.1. Normative References

- [I-D.ietf-rift-rift]  
Przygienda, T., Sharma, A., Thubert, P., Rijsman, B., and D. Afanasiev, "RIFT: Routing in Fat Trees", 2021,  
<<https://www.ietf.org/internet-drafts/draft-ietf-rift-rift.txt>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119,  
DOI 10.17487/RFC2119, March 1997,  
<<https://www.rfc-editor.org/info/rfc2119>>.

### 6.2. Informative References



[RFC7991] Hoffman, P., "The "xml2rfc" Version 3 Vocabulary",  
RFC 7991, DOI 10.17487/RFC7991, December 2016,  
<<https://www.rfc-editor.org/info/rfc7991>>.

Authors' Addresses

Zheng Zhang  
ZTE Corporation  
China

Email: [zhang.zheng@zte.com.cn](mailto:zhang.zheng@zte.com.cn)

Yuehua Wei  
ZTE Corporation  
China

Email: [wei.yuehua@zte.com.cn](mailto:wei.yuehua@zte.com.cn)

Benchong Xu  
ZTE Corporation  
China

Email: [xu.benchong@zte.com.cn](mailto:xu.benchong@zte.com.cn)