# Routing in Fat Trees
# Engineering Simplicity in IP Fabric

Melchior Aelmans

melchior@juniper.net

# AGENDA

- IP Fabrics: The Big Picture
  - Standardized Form-Factor of Network Capacity
  - Next-Gen Evolution Drivers

- Next Gen Underlay Routing Protocol Requirements

- RIFT Basic Concept

- RIFT Advantages

- Standards Status

- Product Status

# THE BIG PICTURE

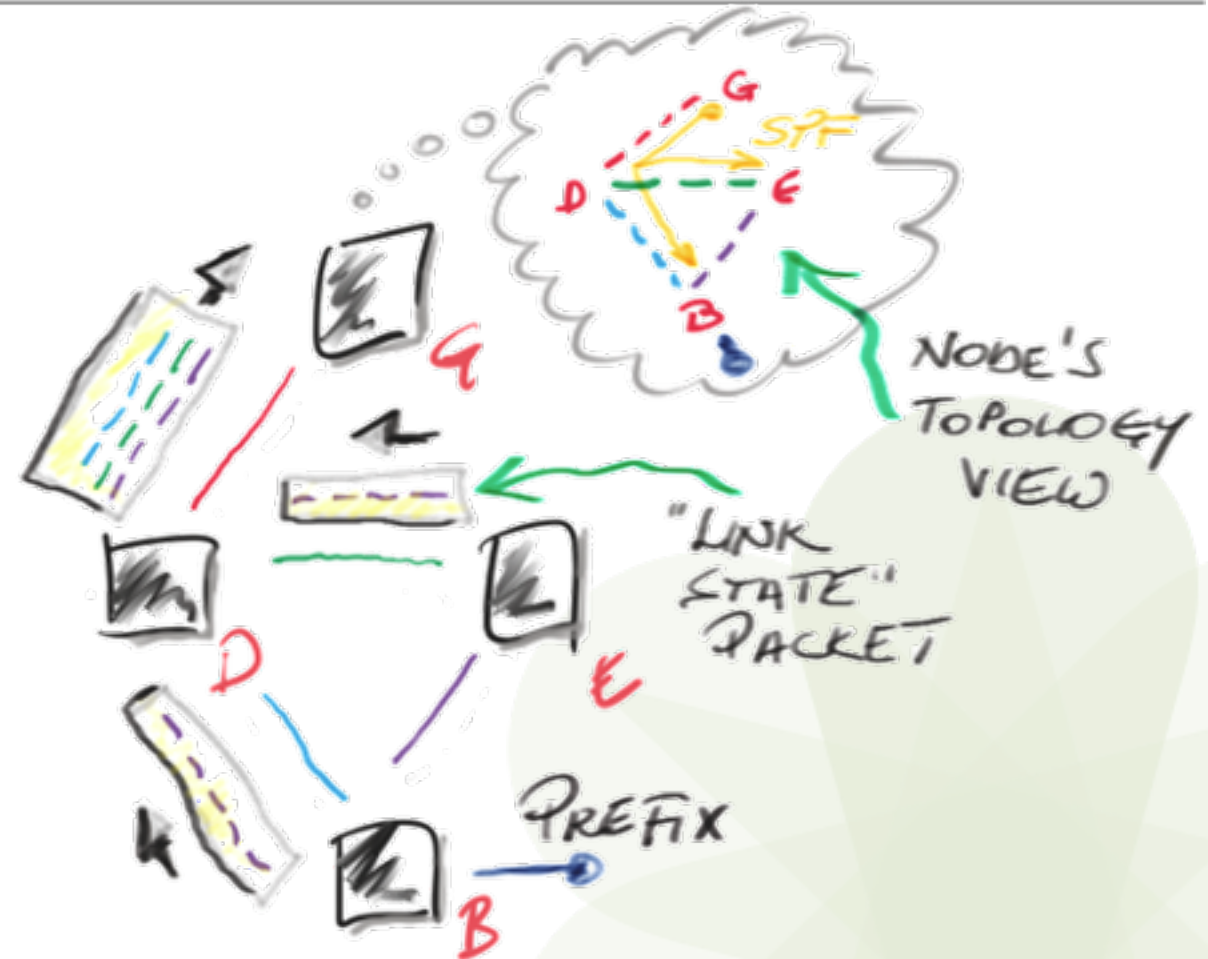## NEXT GEN IP FABRIC EVOLUTION DRIVERS

# IP Fabrics: The Big Picture

- The Global Data Center Market Size is Expected to Reach Revenues of around $174 Billion by 2023
  - growing at a CAGR of about 4% 2018–2023

- Topologies & Technologies are Unifying Towards a "IP Fabric Substrate" in Proximity to Producers and Consumers of Content

- Fortune 1000 Generate and Host Content and are Building Own IP Fabrics to Consume Larger and Larger Amounts of "Local" Networking Capacity

- Building IP Fabrics Today is Largely High OPEX Artisanal Activity

- Most Fortune 1000 Need a "Standardized Form Factor" for Bandwidth Just Like Storage

- We Lack a Standardized ZTP Underlay Technology to Make IP Forwarding in Fabrics as Simple as Ethernet Switching

# WHAT AND WHY ?

- Hyper-scalers are Extrapolating the Things to Come
  - Vast Amount of Bandwidth Close to Producer & Consumer Are Provisioned
    - IP Fabrics in DC (Server Farms)
    - Metro (Caches and Access)
    - Disaggregated Chassis Architectures in Backbone with Regular Fabrics Holding Leaves Together
  - Those Topologies are Uniform, Local and Regular
    - At the End, Economics of Interconnecting of Crossbars Done in 1950 in Bell Labs are Still Valid
  - WAN-Style Traffic Engineering & Protection is Being Replaced by Wide Fan-Out & Distributed Systems Redundancy (Rather Than Chassis & FRR)
    - Simpler is Cheaper in Opex _and_ Capex
  - Hyper-Scalers are Building Customized High-Opex Solutions to Manage those Fabrics

- IP Fabric is Becoming the New "RAM Chip" to Consume Bandwidth
  - No'one Configures SSD Wear-Leveling, RAM Banks and CAS/RAS Manually in Every Laptop
  - IP Fabrics HW is Largely Commodity Already
    - L3 Forwarding is the New L2
  - IP Fabrics Must and Will "OPEX Commoditize"
- Customers are Hosting Their Content & Critical Business Processes
  - Hybrid Cloud for Many Reasons, One of Them to Keep Real-Estate from Hyper-scalers
  - Need to Build Own Fabrics
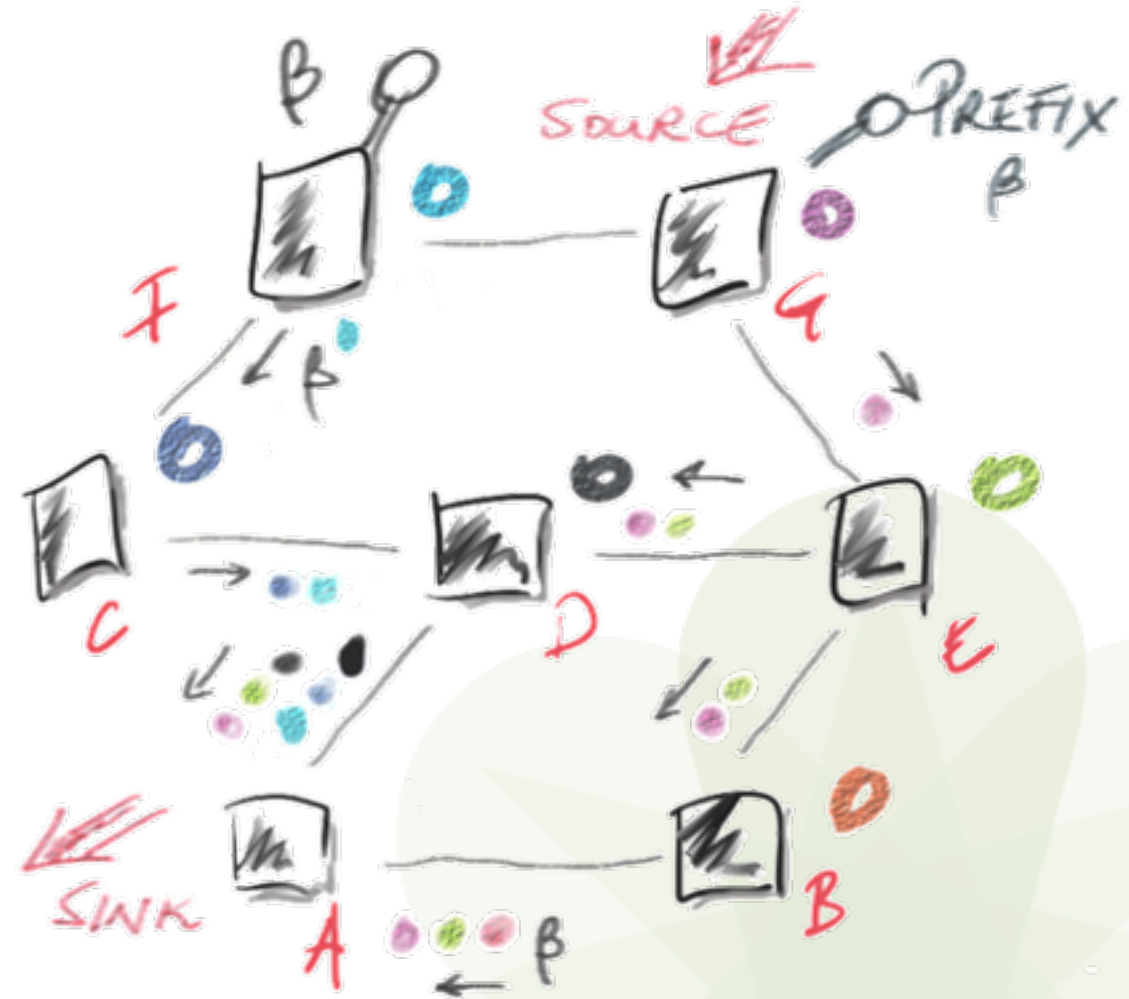  - Hard to Sustain Proprietary OPEX Efforts

# LINK STATE AND SPF = *DISTRIBUTED* COMPUTATION

- TOPOLOGY ELEMENTS
  - NODES
  - LINKS
  - PREFIXES
- EACH NODE ORIGINATES PACKETS WITH ITS ELEMENTS
- PACKETS ARE "FLOODED"
- "NEWEST" VERSION WINS
- EACH NODE "SEES" WHOLE TOPOLOGY
- EACH NODE "COMPUTES" REACHABILITY TO EVERYWHERE
- CONVERSION IS VERY FAST
- EVERY LINK FAILURE SHAKES WHOLE NETWORK (MODULO AREAS)
- FLOODING GENERATES EXCESSIVE LOAD FOR LARGE AVERAGE CONNECTIVITY
- PERIODIC REFRESHES (NOT STRICTLY NECESSARY)

# DISTANCE/PATH VECTOR = *DIFFUSED* COMPUTATION

- PREFIXES "GATHER" METRIC WHEN PASSED ALONG LINKS

- EACH SINK COMPUTES "BEST" RESULT AND PASSES IT ON ( ADD-PATH CHANGED THAT )

- A SINK KEEPS ALL COPIES, OTHERWISE IT WOULD HAVE TO TRIGGER "RE-DIFFUSION"

- LOOP PREVENTION IS EASY ON STRICTLY UNIFORMLY INCREASING METRIC

- IDEAL FOR "POLICY" RATHER THAN "REACHABILITY"

- SCALES WHEN PROPERLY IMPLEMENTED TO MUCH HIGHER # OF ROUTES THAN LINK-STATE

# CURRENT STATE OF AFFAIRS

- SEVERAL OF LARGE DC FABRICS USE E-BGP WITH BAND-AIDS AS DE-FACTO IGP (RFC7938)
  - NUMBERING SCHEMES TO CONTROL "PATH HUNTING"
    - "LOOPING PATHS" (ALLOW-OWN-AS UNDER AS PRIVATE NUMBERING)
    - "RELAXED MULTI-PATH ECMP" SINCE ECMP OVER DIFFERENT AS IN EBGP DOES NOT WORK NORMALLY
  - ADD PATHS TO SUPPORT MULTI-HOMING, N-ECMP, PREVENT OSCILLATIONS
  - EFFORTS TO GET AROUND 65K ASES AND LIMITED PRIVATE AS SPACE
  - PROPRIETARY PROVISIONING AND CONFIGURATION SOLUTIONS, LLDP EXTENSIONS
  - "VIOLATIONS" OF FSM LIKE RESTART TIMERS AND MINIMUM-ROUTE-ADVERTISEMENT TIMERS
  - EMERGING WORK FOR "PEER AUTO-DISCOVERY" AND "SPF" DIAMETRICALLY OPPOSITE TO BGP DESIGN PRINCIPLES
  - RELIANCE ON "UPDATE GROUPS" ~ PEER GROUPS TO PREVENT WITHDRAWAL AND PATH HUNTING AFTER SERVER LINK FAILURES

- OTHERS RUN IGP (ISIS)
  - GENERALLY A "BETTER" APPROACH TO FASTER CONVERGENCE
  - CURRENT ATTEMPTS TO DEAL WITH SOME "SPOT PROBLEMS" LIKE FLOODING REDUCTION

- YET OTHERS RUN BGP OVER IGP (TRADITIONAL ROUTING ARCHITECTURE)

- LESS THAN MORE SUCCESSFUL ATTEMPTS @ PREFIX SUMMARIZATION, MICRO- AND BLACK-HOLING, BLAST RADIUS CONTAINMENT

- SERVER MULTI-HOMING NOT POSSIBLE USING IP DUE TO EQUAL COST AND SCALING CONSTRAINTS, HENCE MC-LAG'ED SOLUTIONS OR EVPN

- IN SUMMARY: HIGH OPEX SOLUTIONS NOT NECESSARILY VIABLE FOR CUSTOMERS WHO CANNOT OR DO NOT WANT TO BUILD SOPHISTICATED TALENT POOL TO DEAL WITH THEIR "UNICORN" FABRICS

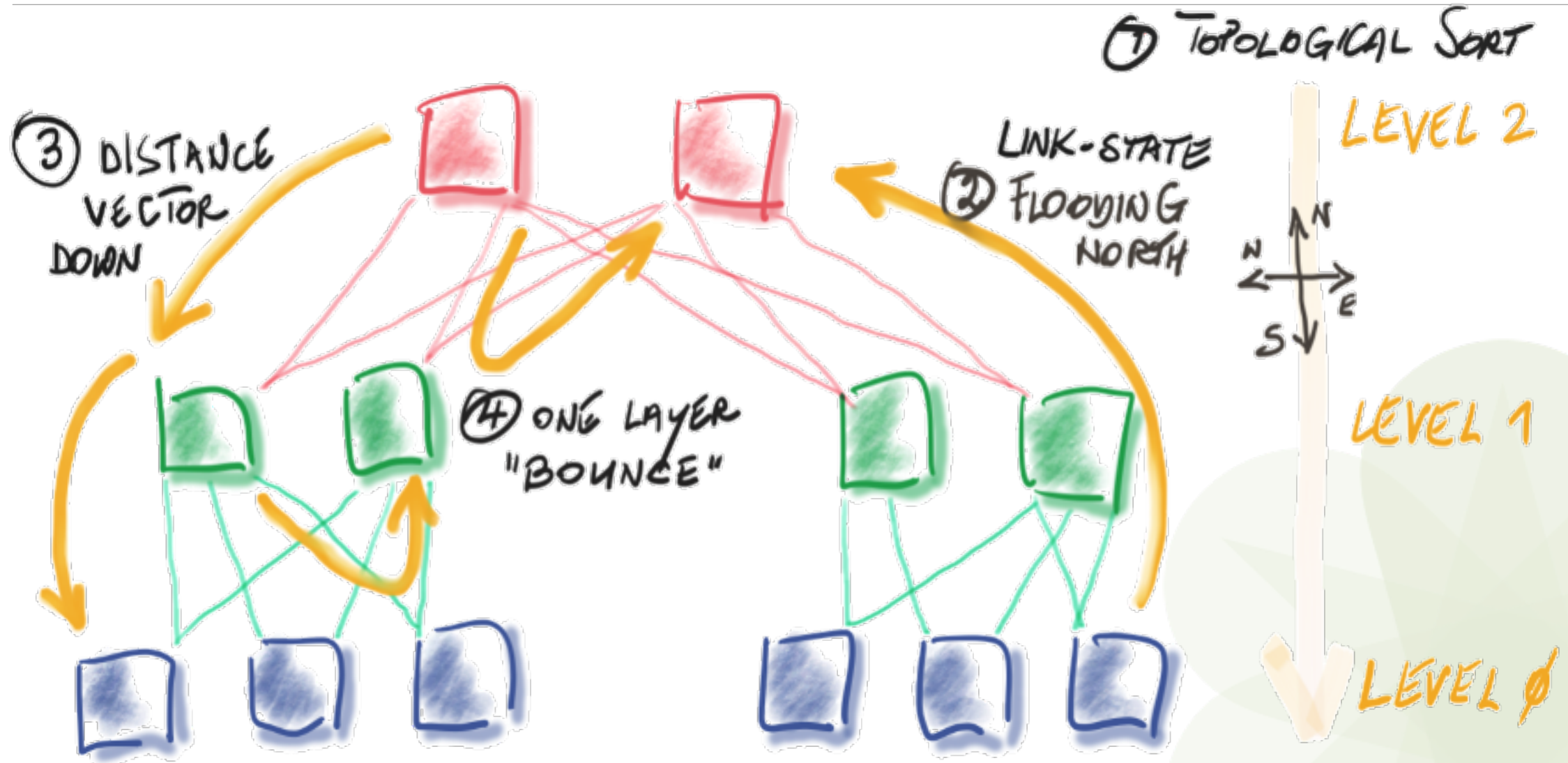# Next Gen IP Fabric Underlay Routing Protocol Requirements

# NEXT GEN UNDERLAY PROTOCOL REQUIREMENTS

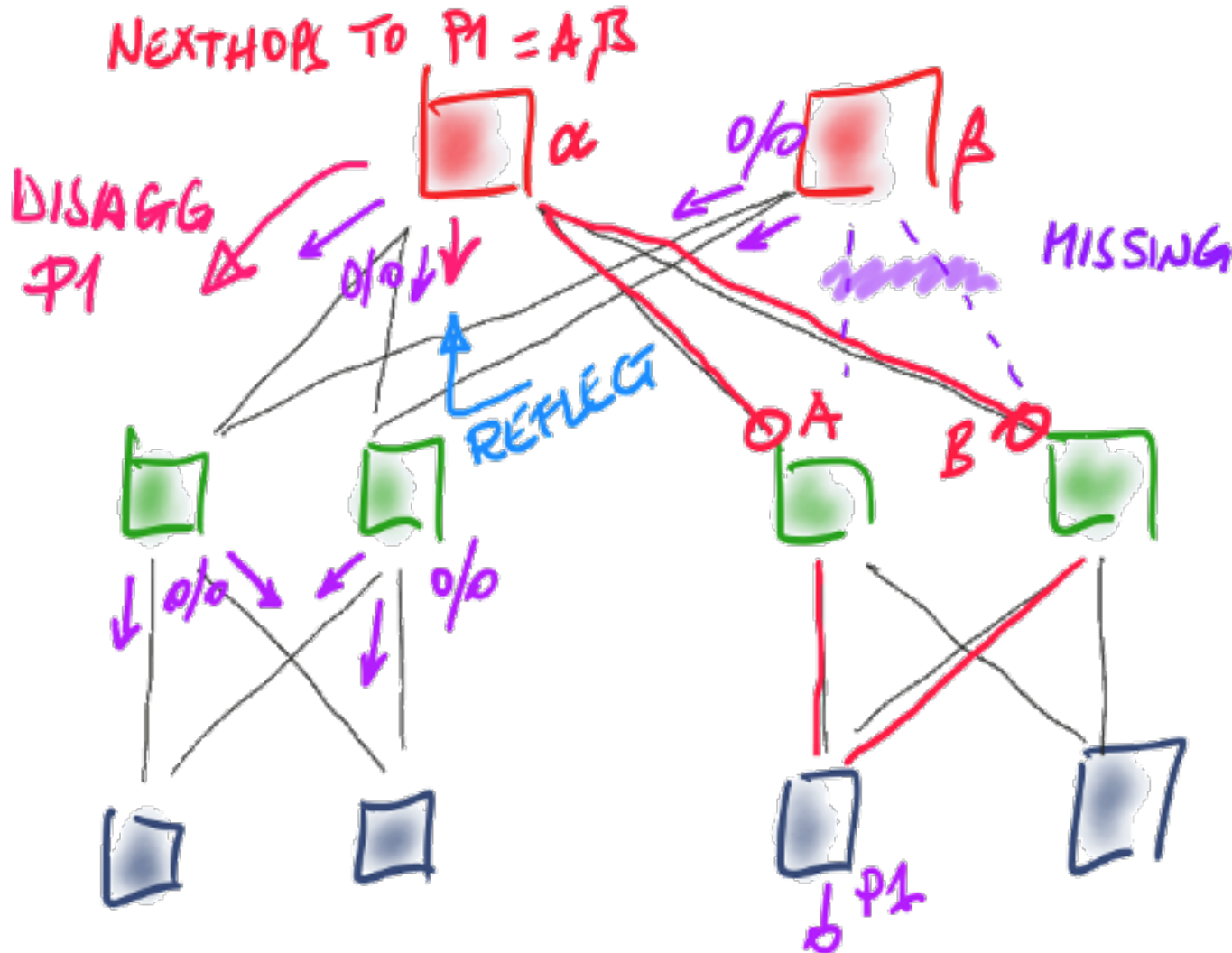| Problem / Attempted Solution | BGP modified for DC (all kind of "mods") | ISIS modified for DC (RFC7356 + "mods") | RIFT Native DC |
|---|---|---|---|
| Peer Discovery/Automatic Forming of Trees/Preventing Cabling Violations | ⚠️ | ⚠️ | ✓ |
| No Need for Internal Addressing, Minimal Amount of Routes/Information on ToRs, Stretches True Routing to the Multi-Homed Host | ✗ | ✗ | ✓ |
| High Degree of ECMP (BGP needs lots knobs, memory, own-AS-path violations) and ideally NEC and LFA | ⚠️ | ✓ | ✓ |
| Non Equal Cost Multi-Path, ECMP Independent Anycast, MC-LAG Replacement | ✗ | ✗ | ✓ |
| Traffic Engineering by Next-Hops, Prefix Modifications | ✓ | ✗ | ✓ |
| See All Links in Topology to Support PCE/SR | ⚠️ | ✓ | ✓ |
| Carry Opaque Configuration Data (Key-Value) Efficiently | ✗ | ⚠️ | ✓ |
| Take a Node out of Production Quickly and Without Disruption | ✗ | ✓ | ✓ |
| Automatic Disaggregation on Failures to Prevent Black-Holing and Back-Hauling | ✗ | ✗ | ✓ |
| Minimal Blast Radius on Failures (On Failure Smallest Possible Part of the Network "Shakes") | ✗ | ✗ | ✓ |
| Fastest Possible Convergence on Failures | ✗ | ✓ | ✓ |
| State of the Art Security Including Originator Validation and Replay Prevention | ✗ | ✗ | ✓ |
| Simplest Initial Implementation | ✓ | ✗ | ✓✗ |

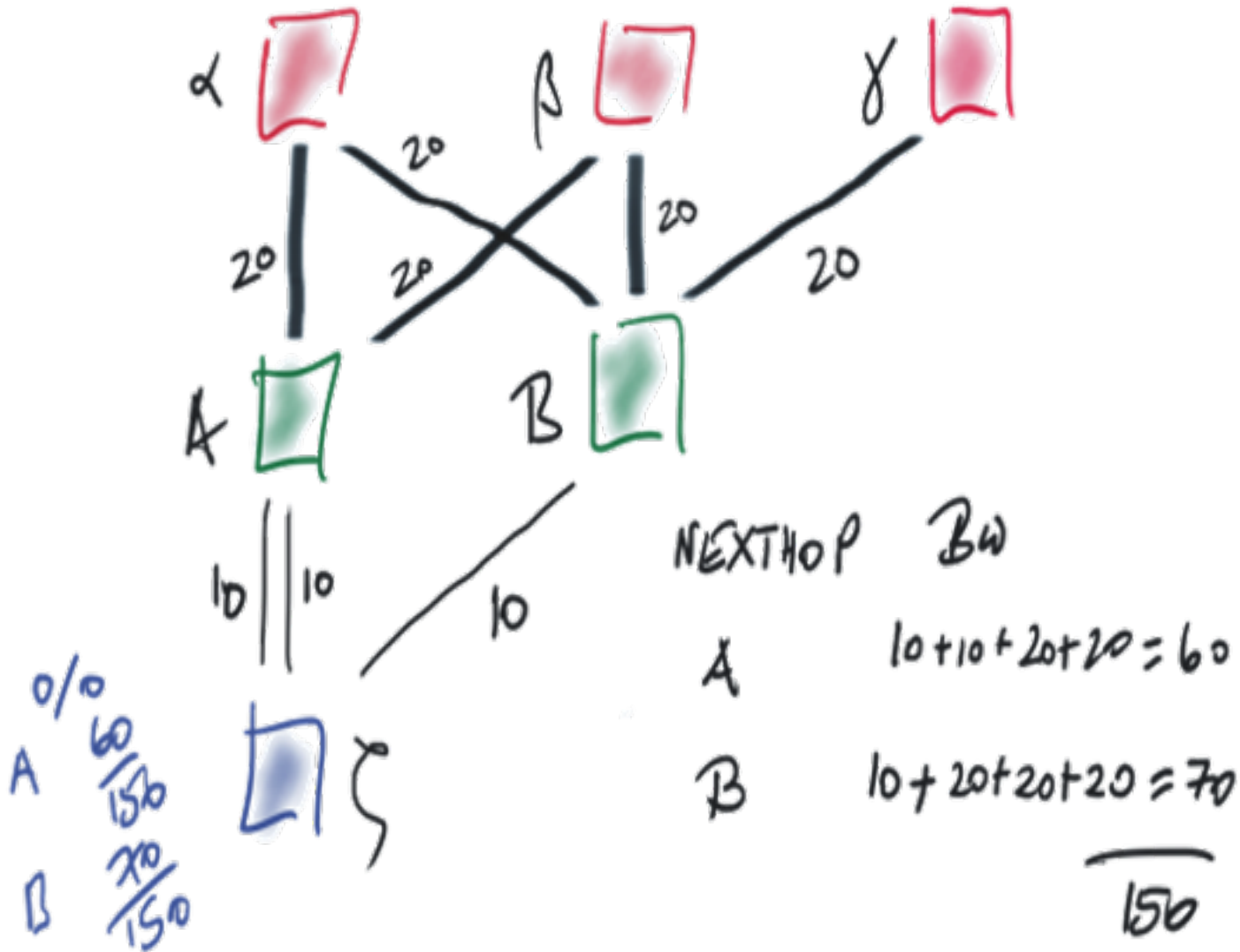JUNIPER NETWORKS

# RIFT BASIC CONCEPT

# RIFT Concept

# RIFT Advantages

# Automatic De-Aggregation



- South Representation with Adjacencies of Red Beta is Reflected by the Green Layer to Alpha
- Alpha Red Spine Computes P1 Reachability (Tree that Indicates Next-Hops)
- alpha Sees that Beta has no Adjacencies to any of the available Next-Hops to P1
- Alpha Disaggregates P1
- Left Green Nodes do NOT Propagate Disaggregation Further
- When Beta Obtains a Viable Next-Hop to P1, Alpha Automatically Reaggregates

# Northbound Bandwidth Balancing



- Zeta Computes Bandwidth Sum to NextHop and Next Level Up
- Through A Cumulative 60 Units
- Through B Cumulative 70 Units
- Default Route Uses That in Some Form to Unequal Load Balance
- Perfectly Safe Since RIFT is Loop-Free

# OPERATIONAL ADVANTAGES

- Open IETF Standard
  - Can Build Hybrid Vendor Fabrics
  - Protocol is Well Reviewed and Understood by World-Class Experts

- True ZTP
  - No Configuration Necessary
  - V4 over V6 Forwarding
  - Mis-cabling Handled

- Can Operate on Asymmetric Bandwidth Fabrics and Handle "Fat Link" Failures By Adjusting Automatically

- Can Support and Scale to an Architecture with Multi-Home Servers
  - No Need for Service Migration on ToR Upgrades
  - Can Talk Directly to Hyper-Visors/Kubernetes GW

- BFD is "Built In"
  - Can Be Used for Fast Rehash or Early Loss Detection

- Runs on UDP
  - Trivial Kernel Support on All Platforms
  - Allows for Max. Speed Flooding
  - Easy to "Multi-Instantiate" for Different Purposes

- Minimal Blast-Radius
  - Failures/Bring-Up on Fabric Only Affects the Smallest Viable Radius

- RIFT Flooding is ~30% of Normal Flat IGP
  - Built-In Flood Reduction Reduces Flood Traffic to <20% of Flat IGP

- Loop-Free
  - Can Utilize **All** Viable Paths Through Fabric
  - Can Support True Anycast

- Model Based
  - Much Less Possibility for Parser and Formatter Bugs Plaguing Today's Networking Protocols

- Specification is Written for Maximum parallelization
  - With Enough Cores IP Switches Should Be Able to Converge @ Speeds Making FRR Unnecessary (Assuming Fast Rehash)

- KV Store Allows to Replace Out-Of-Band Applications
  - IP/MAC Binding Can Be Flooded to Top-of-Fabric

- Sophisticated, Newest Routing Security

- Full Fabric Visibility at the Top

# STANDARDIZATION PROGRESS

# STANDARDIZATION PROGRESS

RIFT STANDARDS TRACK IETF WORKING GROUP IN REVISION -12

- SECURITY AREA DIRECTORATE EARLY REVISION DONE AND ADDRESSED
- IESG REVISIONS BEING ADDRESSED
- CO-AUTHORS FROM CISCO, YANDEX, MELLANOX, HP & IMPLEMENTING INDIVIDUALS
- HTTPS://TOOLS.IETF.ORG/HTML/DRAFT-IETF-RIFT-RIFT-12

APPLICABILITY DRAFT

- CO-AUTHORED WITH OPERATORS, E.G. ORANGE AND YANDEX
- UNDER GEN-ART REVIEW
- HTTPS://TOOLS.IETF.ORG/HTML/DRAFT-IETF-RIFT-APPLICABILITY-03

# PRODUCT STATUS

# PRODUCTIZATION

**IMPLEMENTATIONS**

**JUNIPER NETWORKS**

- FCS IN 19.4R1 64-BIT JUNOS
- PACKAGE INSTALLING OVER 19.4R1 OR NEWER
- QFX (ALL VARIANTS), vMX, MX
- cRPD & EVO IN PROGRESS
- https://www.juniper.net/us/en/dm/free-rift-trial/

**PYTHON OPEN-SOURCE BY BRUNO RIJSMAN**

- https://github.com/brunorijsman/rift-python

**DOCUMENTATION**

- DAY ONE BOOK: ROUTING IN FAT TREES (RIFT)
  - COVERING JUNIPER AND OPEN-SOURCE IMPLEMENTATION
  - RIFT WIRESHARK DISSECTOR
  - DOWNLOAD FREE: HTTP://JNPR.NL/RIFT

JUNIPER | Engineering
NETWORKS | Simplicity

DAY ONE: ROUTING IN FAT TREES (RIFT)

A complete look at the cutting edge protocol.

By Melchior Aelmans, Olivier Vandezande, Bruno Rijsman,
Jordan Head, Christian Graf, Leonardo Alberro,
Hitesh Mali, Oliver Steudler

# RIFT: Engineering Simplicity in IP Fabric Questions?

Melchior Aelmans

melchior@juniper.net

JUNIPER NETWORKS | Engineering Simplicity