

RIFT: OPEN STANDARD, ZERO OPEX, IP FABRIC ROUTING UNDERLAY

IEEE 802/IETF DC WORKSHOP

TONY PRZYGIENDA
DISTINGUISHED ENGINEER, JUNIPER NETWORKS

DRAFT-IETF-RIFT-RIFT @ IETF

DISCLAIMERS AND EXPECTATIONS

- THIS IS AN IETF “WORKING STRAW-MAN PROPOSAL”
- UNLESS SPECIFICALLY STATED NONE OF THOSE THINGS CONSTITUTE COMMITMENTS TO PRODUCT SPECIFICATIONS, OFFERINGS OR RELEASE DATES BY JUNIPER NETWORKS AT THIS POINT IN TIME

WHAT AND WHY ?

- HYPER-SCALERS ARE EXTRAPOLATING THE THINGS TO COME
 - VAST AMOUNT OF BANDWIDTH CLOSE TO PRODUCER & CONSUMER NECESSARY
 - IP FABRICS IN DC (SERVER FARMS)
 - METRO (CACHES AND ACCESS)
 - DISAGGREGATED CHASSIS ARCHITECTURES
 - THOSE TOPOLOGIES ARE BECOMING UNIFORM, LOCAL AND REGULAR
 - WAN-STYLE TRAFFIC ENGINEERING & PROTECTION IS BEING REPLACED BY WIDE FAN-OUT & DISTRIBUTED SYSTEMS REDUNDANCY (RATHER THAN CHASSIS & FRR)
 - 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 RAM BANKS AND CAS/RAS MANUALLY IN EVERY LAPTOP
 - IP FABRICS HW IS LARGELY COMMODITY ALREADY
 - IP FABRICS MUST “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

AGENDA

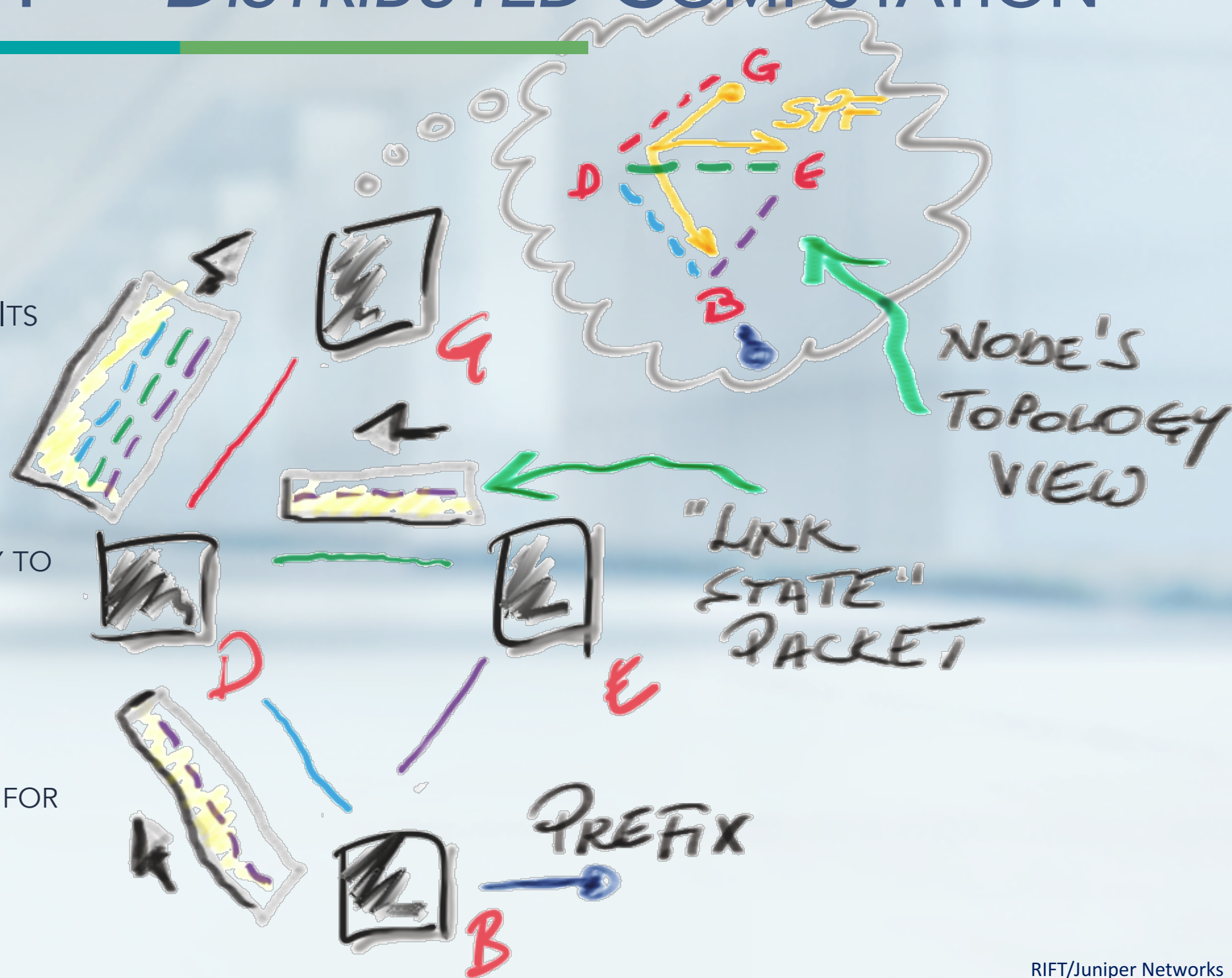
- BLITZ OVERVIEW OF TODAY'S ROUTING (IF NEEDED)
- "FABRIC ROUTING" IS A SPECIALIZED PROBLEM
- RIFT: A NOVEL ROUTING ALGORITHM FOR IP FABRIC UNDERLAY

BLITZ OVERVIEW OF TODAY'S ROUTING

- LINK STATE & SPF
- DISTANCE/PATH VECTOR

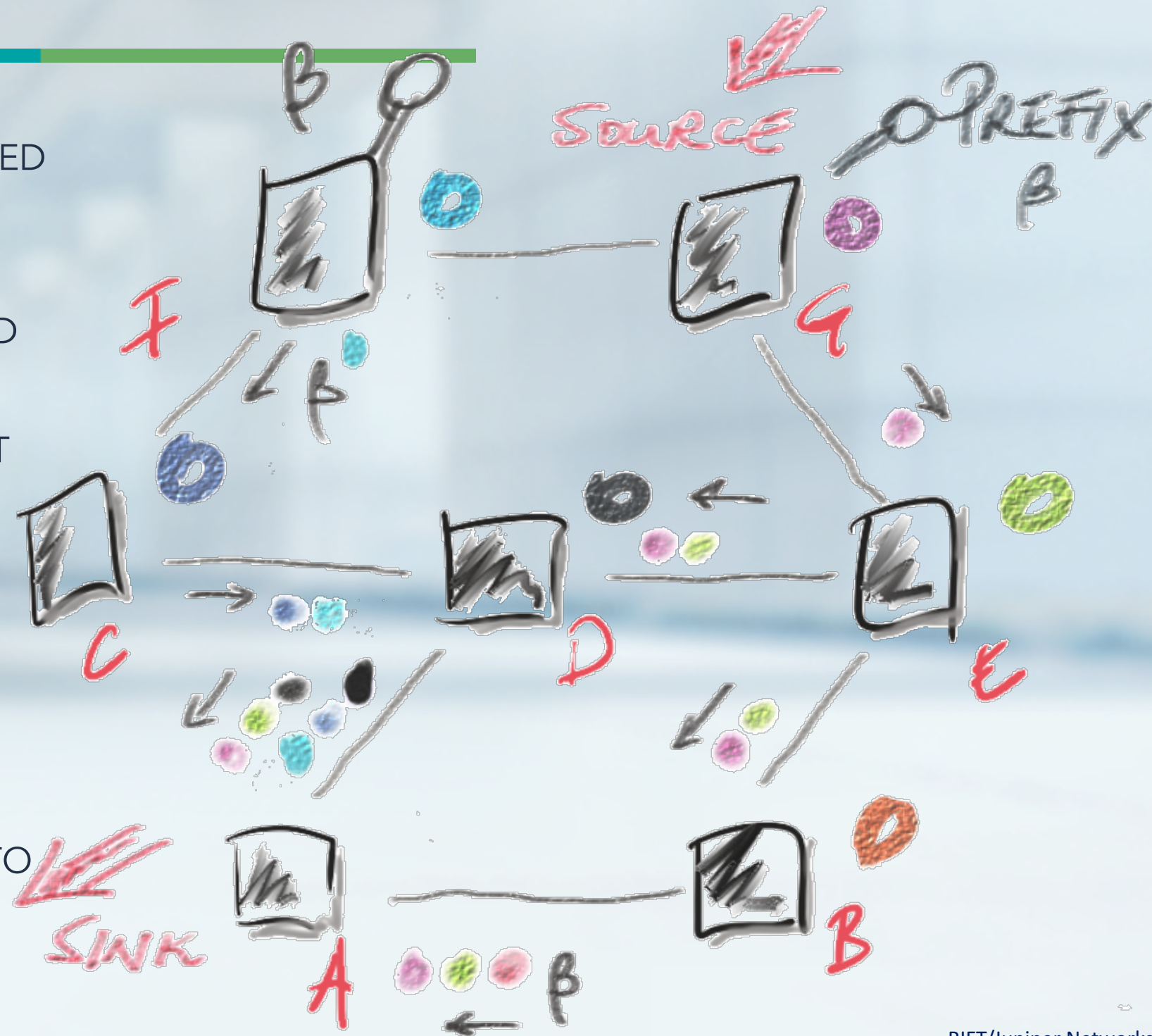
LINK STATE AND SPF = *DISTRIBUTED* COMPUTATION

- 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

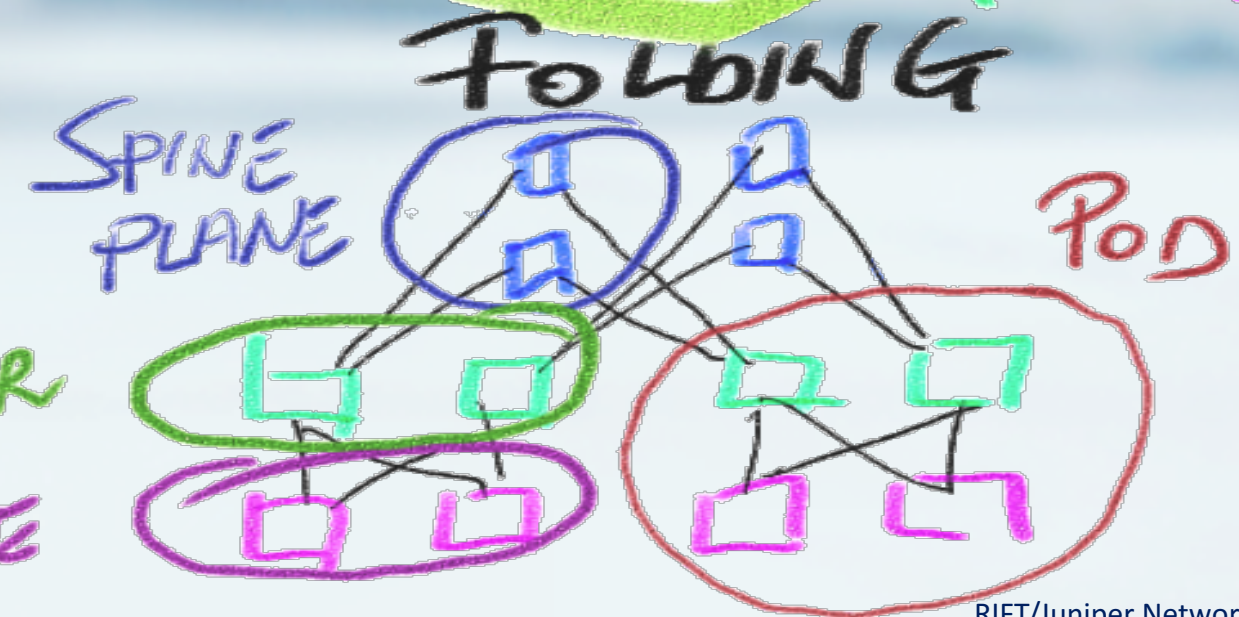
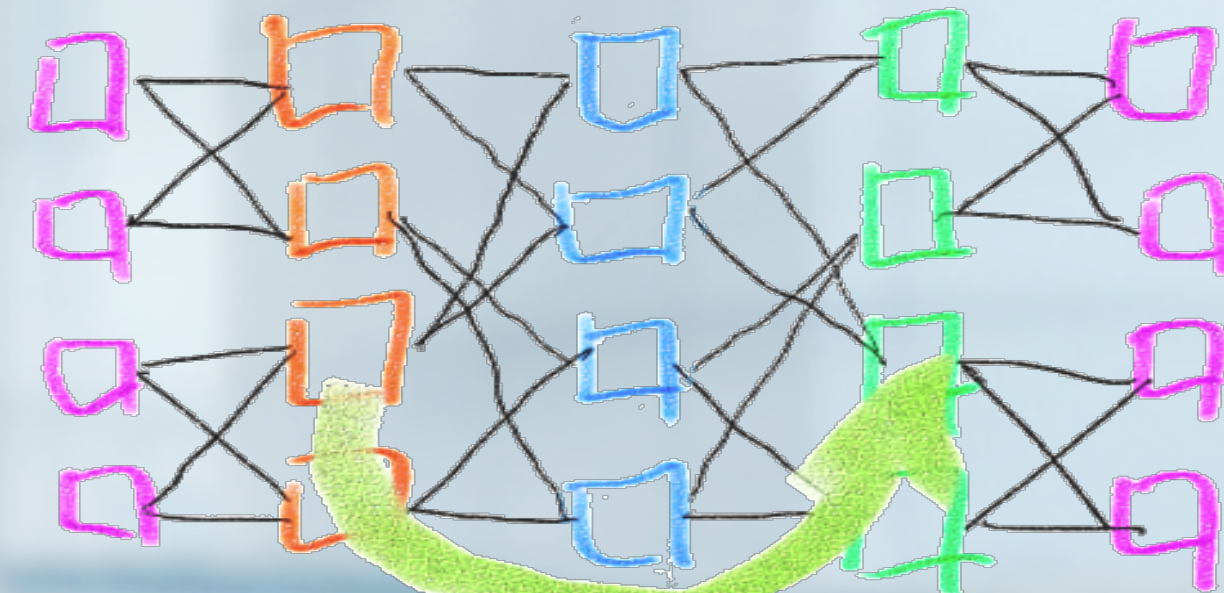
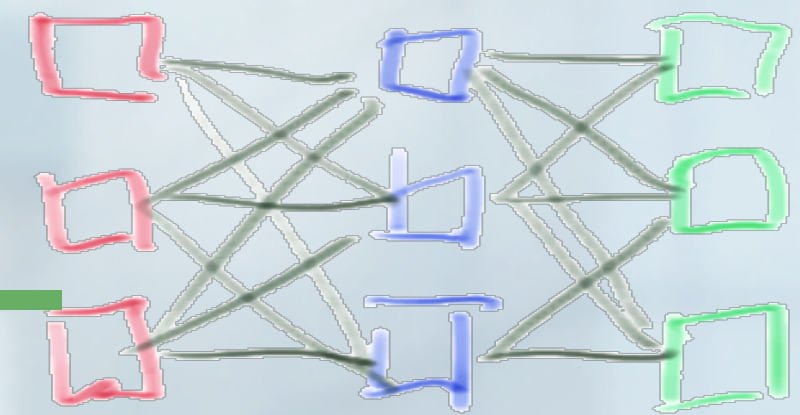


DC FABRIC ROUTING: A SPECIALIZED PROBLEM

- CLOS TOPOLOGIES ARE DOMINANT TODAY
 - TOROIDAL [AND DIAGONAL] MESHES HAVE LONG PATHS, SMALL BISECTION WIDTH AND POOR BLOCKING PROPERTIES
 - DRAGONFLY (AND SOME PROBABILISTIC VARIANTS) IS VERY NOVEL AND UNPROVEN
 - $\frac{1}{2}$ THROUGHPUT OF CLOS AT SAME COST DUE TO LOW ECMP
 - RIFT COULD WORK WELL IN A PRACTICAL MODIFICATION (ONE LEVEL CLOS AND DRAGONFLY CORE)
- CURRENT STATE OF AFFAIRS
- REQUIREMENTS MATRIX

CLOS TOPOLOGIES

- CLOS OFFERS WELL-UNDERSTOOD BLOCKING PROBABILITIES
- WORK DONE AT AT&T (BELL SYSTEMS) IN 1950s
- FULLY CONNECTED CLOS IS DENSE AND EXPENSIVE
- DATA CENTERS TODAY TEND TO BE VARIATIONS OF "FOLDED FAT-TREE"
 - INPUT STAGES = OUTPUT STAGES
 - CLOS IS "PARTIAL"
 - LINKS GET "FATTER" UP THE TREE



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

REQUIREMENTS BREAKDOWN (RFC7938+) FOR A "MINIMAL OPEX FABRIC"

Problem / Attempted Solution	BGP modified for DC (all kind of "mods")	ISIS modified for DC (RFC7356 + "mods")	RIFT Native DC
Optional Peer Discovery/True ZTP/Preventing Cabling Violations	⚠	⚠	✓
Minimal Amount of Routes/Information on ToRs, light-weight enough for servers, Can Scale to Multi-Homed Server Architectures	✗	✗	✓
High Degree of ECMP (BGP needs lots knobs, memory, own-AS-path violations) and ideally NEC and LFA	⚠	✓	✓
Non Equal Cost Multi-Path, Equal Cost 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	✗	✓	✓
Bandwidth Load Balancing	✗	✗	✓
Simplest Initial Implementation	✓	✗	✗

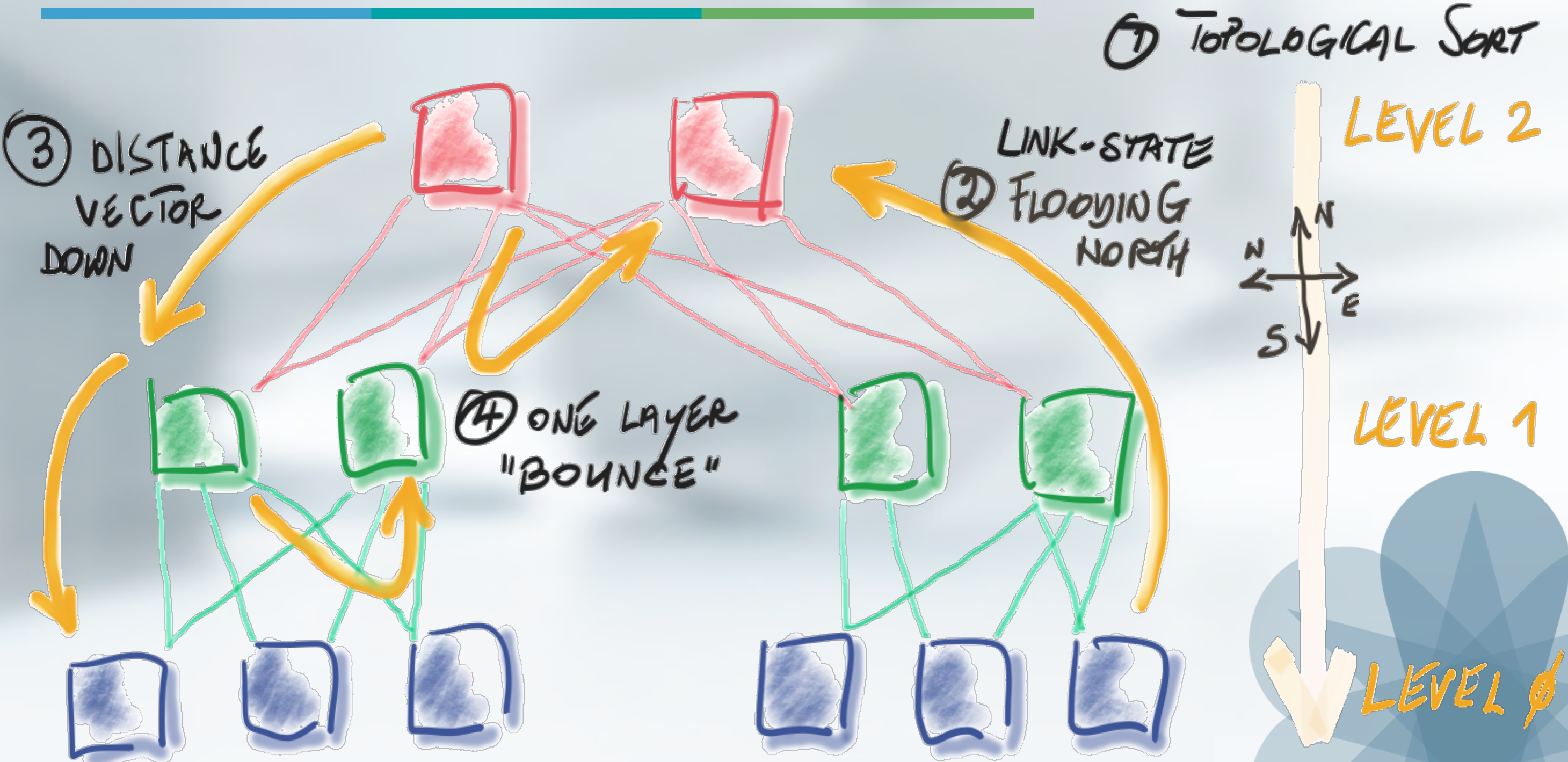
RIFT: NOVEL ROUTING ALGORITHM FOR CLOS UNDERLAY

- GENERAL CONCEPT
- AUTOMATIC DISAGGREGATION
- AUTOMATIC BANDWIDTH BALANCING
- FAST MOBILITY SUPPORT
- AND MORE

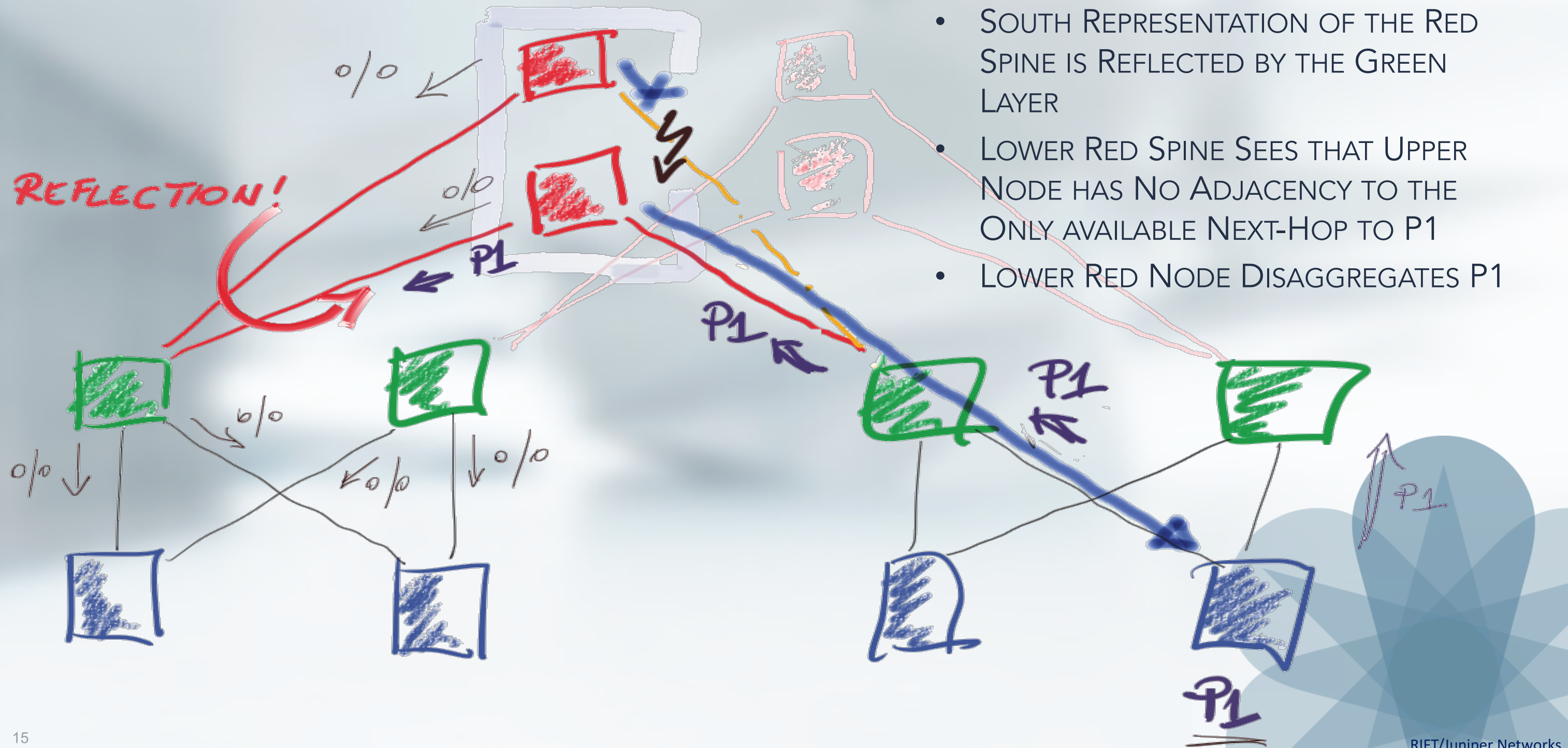
“Just because the standard provides a cliff in front of you, you are not necessarily required to jump off it.”

— Norman Diamond

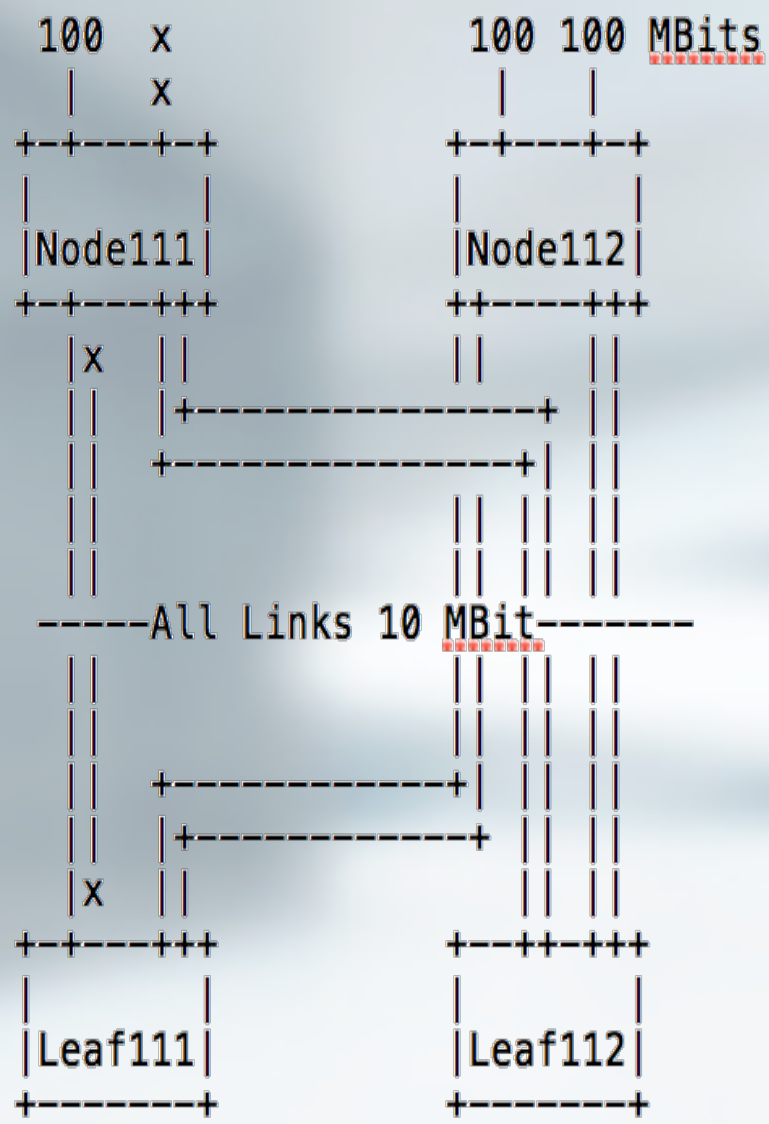
LINK-STATE UP, DISTANCE VECTOR DOWN & BOUNCE



AUTOMATIC DE-AGGREGATION



NORTHBOUND BANDWIDTH BALANCING



- RIFT calculates the amount of northbound bandwidth available towards a node compared to other nodes at the same level and adjusts the default route distance accordingly to allow for the lower level to have different weights on load balancing.
- **BAD_N**: Bandwidth Adjusted Metric to N
- **L_N_u**: as sum of the bandwidth available from L to N
- **N_u**: as sum of the uplink bandwidth available on N
- **T_N_u**: $L_N_u + N_u$
- **M_N_u**: $\log_2(\text{next_power_2}(T_N_u))$
- **BAD_N**: $D * (1 + \text{maximum_of_all}(M_N_u) - M_N_u)$

Node	N	T_N_u	M_N_u	BAD
Leaf111	Node111	110	7	2
Leaf111	Node112	220	8	1
Leaf112	Node111	120	7	2
Leaf112	Node112	220	8	1

MOBILITY SUPPORT

- OPTIONAL CLOCK ATTRIBUTE ON PREFIX
- IF CLOCK NOT PRESENT, ADDRESS IS ANYCAST
- IF PRESENT, ALWAYS BETTER THAN NONE
- IF BOTH PRESENT, RFC5905 OR BETTER ON FABRIC ASSUMED
 - IF IEEE802_1 LESS THAN 200MSEC DIFF, TRANSACTIONID (TID) IF PRESENT TIE-BREAKS
 - OTHERWISE TIMESTAMP COMPARES
- TIDS ARE COMING FROM [DRAFT-IETF-6LO-RFC6775-UPDATE](#) OR SIMILAR MECHANISMS

AND THEN ...

- COMPLETE ZTP
 - NOT EVEN ADDRESSING NECESSARY (EXCEPT V6 LOCAL VIA ND)
 - IPV4 OVER IPV6 FORWARDING
 - ARBITRARY NUMBER OF LEVELS
 - HETEROGENEOUS POD HEIGHT POSSIBLE
- LOOP-FREE, I.E. ALL PATHS THROUGH IP FABRIC CAN BE SATURATED
- NORMAL OPERATION HAS ONLY DEFAULT ROUTES ON LEAFS
- MINIMAL BLAST RADIUS
 - AUTOMATIC OPTIMAL FLOODING PRUNING AND LOAD BALANCING ON CHANGES FOR MAXIMUM SCALING
- NORTHBOUND BANDWIDTH BALANCING ON LINK LOSS
- K/V STORE
- COMPLETELY MODEL BASED PACKET FORMATS
- FLOODING OVER UDP FOR FASTEST CONVERGENCE
- POLICY CONTROLLED KEY-VALUE STORE SUPPORT
- POSSIBLE SR SUPPORT

STANDARDIZATION & OPEN SOURCE

- STANDARDS TRACK WORKING GROUP IN IETF
 - JUNIPER & APSTRA CO-CHAIR
- CISCO, COMCAST, YANDEX, MELLANOX, HPE CO-AUTHORS
 - BLOOMBERG, CRITEO & OTHERS ENGAGED
- YANG, SR & OTHER THINGS UNDER WORKS
- OPEN SOURCE IMPLEMENTATION
 - [HTTPS://GITHUB.COM/BRUNORIJSMAN/RIFT-PYTHON](https://github.com/brunorijsman/rift-python)

SUMMARY OF RIFT 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 SCALE TO AND 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 WEIRD 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

SUMMARY OF RIFT PROTOCOL ADVANTAGES

- OPEN IETF STANDARD
- ADVANTAGES OF LINK-STATE AND DISTANCE VECTOR
 - FASTEST POSSIBLE CONVERGENCE
 - AUTOMATIC DETECTION OF TOPOLOGY
 - MINIMAL ROUTES ON TORs
 - HIGH DEGREE OF ECMP
 - FAST DE-COMMISSIONING OF NODES
 - MAXIMUM PROPAGATION SPEED WITH FLEXIBLE # PREFIXES IN AN UPDATE
- VECTOR
 - REDUCED FLOODING
 - AUTOMATIC NEIGHBOR DETECTION
- UNIQUE RIFT ADVANTAGES
 - BANDWIDTH RE-BALANCING
 - AUTOMATIC DISAGGREGATION ON FAILURES
 - KEY-VALUE STORE
 - HORIZONTAL LINKS USED FOR PROTECTION ONLY
 - MINIMAL BLAST RADIUS ON FAILURES
 - CAN UTILIZE ALL PATHS THROUGH FABRIC WITHOUT LOOPING
 - SUPPORTS NON-EQUAL COST MULTIPATH AND CAN REPLACE MC-LAG
 - TRUE ZTP

MORE MATERIAL

- SPECIFICATIONS IN IETF WORKING GROUP
 - [HTTPS://DATATRACKER.IETF.ORG/DOC/DRAFT-IETF-RIFT-RIFT/](https://datatracker.ietf.org/doc/draft-ietf-rift-rift/)
- WALK THROUGH MAJOR CONCEPTS & PACKAGE EXPLANATION (RIFT INTERIM RECORDING)
 - MAY 3, 2018: [HTTPS://WWW.YOUTUBE.COM/WATCH?V=DTXNoCkC7MA](https://www.youtube.com/watch?v=DTXNoCkC7MA)
 - MAY 2, 2018: [HTTPS://WWW.YOUTUBE.COM/WATCH?V=BZtFPTGCSBS](https://www.youtube.com/watch?v=BZtFPTGCSBS)
- JUNIPER'S PUBLIC STAND-ALONE PACKAGE DOWNLOADABLE
 - [HTTPS://WWW.JUNIPER.NET/US/EN/DM/FREE-RIFT-TRIAL/](https://www.juniper.net/us/en/dm/free-rift-trial/)
- OPEN SOURCE IMPLEMENTATION
 - [HTTPS://GITHUB.COM/BRUNORIJSMAN/RIFT-PYTHON](https://github.com/brunorijsman/rift-python)

THANKS
