

RIFT: A NOVEL DC FABRIC ROUTING PROTOCOL

DRAFT-PRZYGIENDA-RIFT

IETF '98

CONTENT

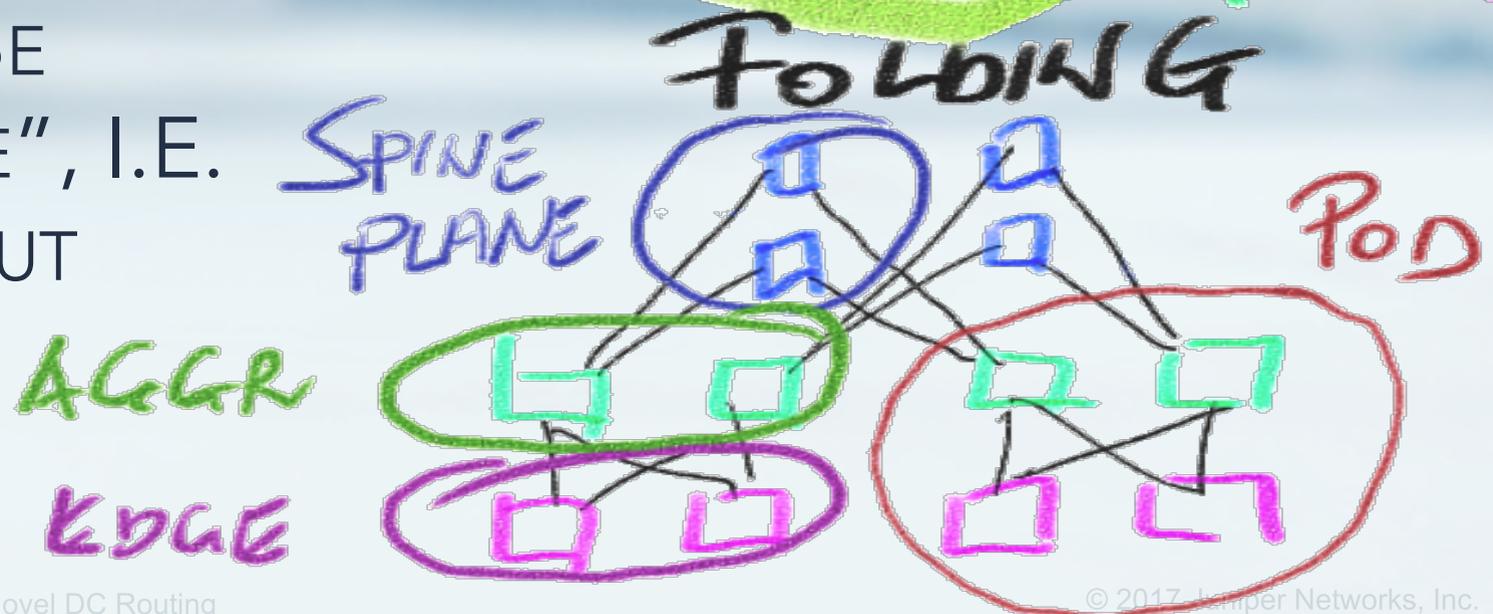
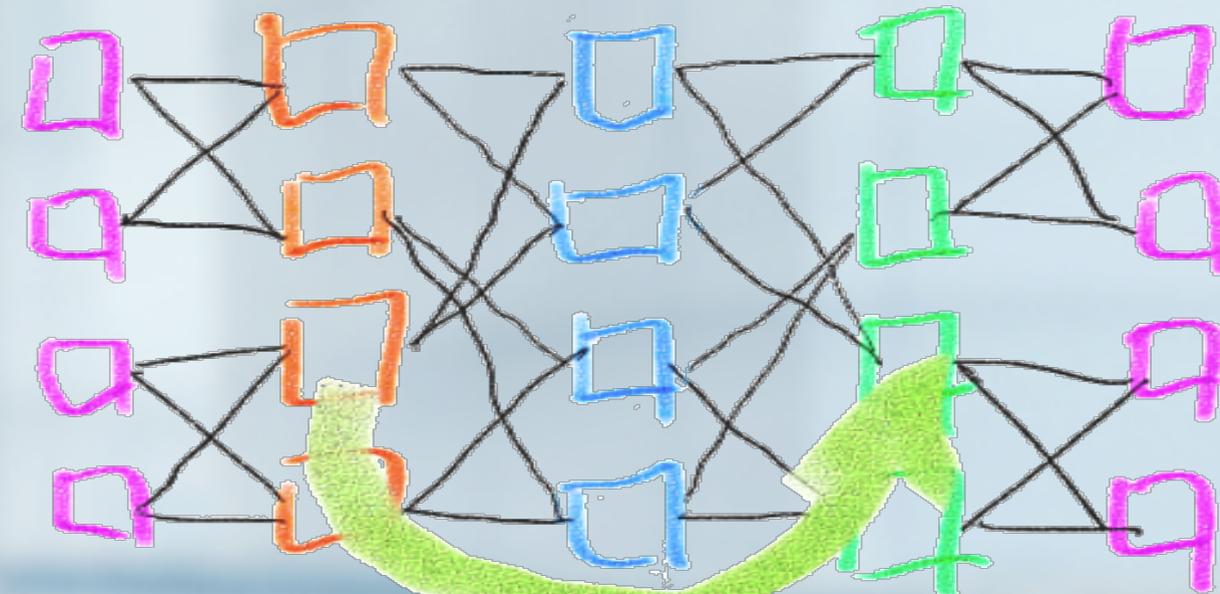
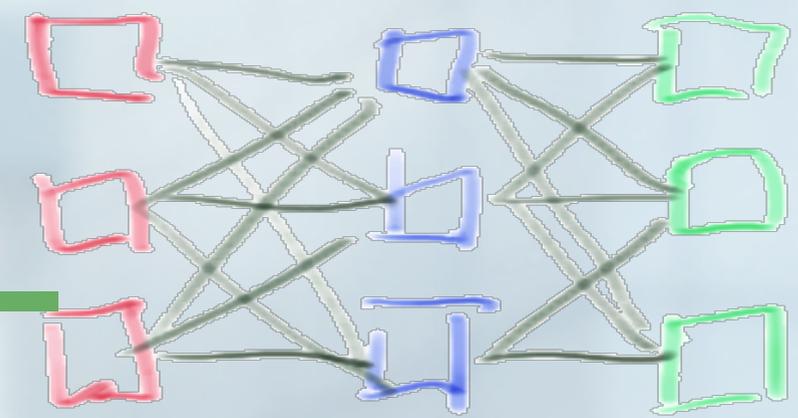
- DC FABRIC ROUTING IS A SPECIALIZED PROBLEM
- RIFT: A NOVEL ROUTING ALGORITHM FOR DC FABRIC UNDERLAY

DC FABRIC ROUTING: A SPECIALIZED PROBLEM

- CLOS/FAT-TREE TOPOLOGY VARIATIONS
- CURRENT STATE OF DYNAMIC DC ROUTING
- DYNAMIC DC ROUTING REQUIREMENTS MATRIX

CLOS VARIATION TOPOLOGIES

- CLOS OFFERS WELL-UNDERSTOOD BLOCKING PROBABILITIES
- WORK DONE AT AT&T (BELL SYSTEMS) IN 1950S FOR CROSSBAR SCALING
- FULLY CONNECTED CLOS IS DENSE AND EXPENSIVE
- DATA CENTERS TODAY TEND TO BE VARIATIONS OF "FOLDED FAT-TREE", I.E. INPUT STAGES ARE SAME AS OUTPUT STAGES AND CLOS IS "PARTIAL"



CURRENT STATE OF AFFAIRS

- SEVERAL OF LARGE DC FABRICS USE E-BGP WITH BAND-AIDS AS IGP (RFC7938)
 - “LOOPING PATHS” (ALLOW-AS)
 - “RELAXED MULTI-PATH ECMP”
 - AS NUMBERING SCHEMES TO CONTROL “PATH HUNTING” VIA POLICIES
 - ADD PATHS TO SUPPORT MULTI-HOMING, ECMP ON EBGP
 - 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
- OTHERS RUN IGP (ISIS)
- YET OTHERS RUN BGP OVER IGP (TRADITIONAL ROUTING ARCHITECTURE)
- LESS THAN MORE SUCCESSFUL ATTEMPTS @ PREFIX SUMMARIZATION, MICRO- AND BLACK-HOLING
 - WORKS BETTER FOR SINGLE-TENANT FABRICS WITHOUT LAN STRETCH OR VM MOBILITY

DYNAMIC DC ROUTING REQUIREMENTS BREAKDOWN (RFC7938+)

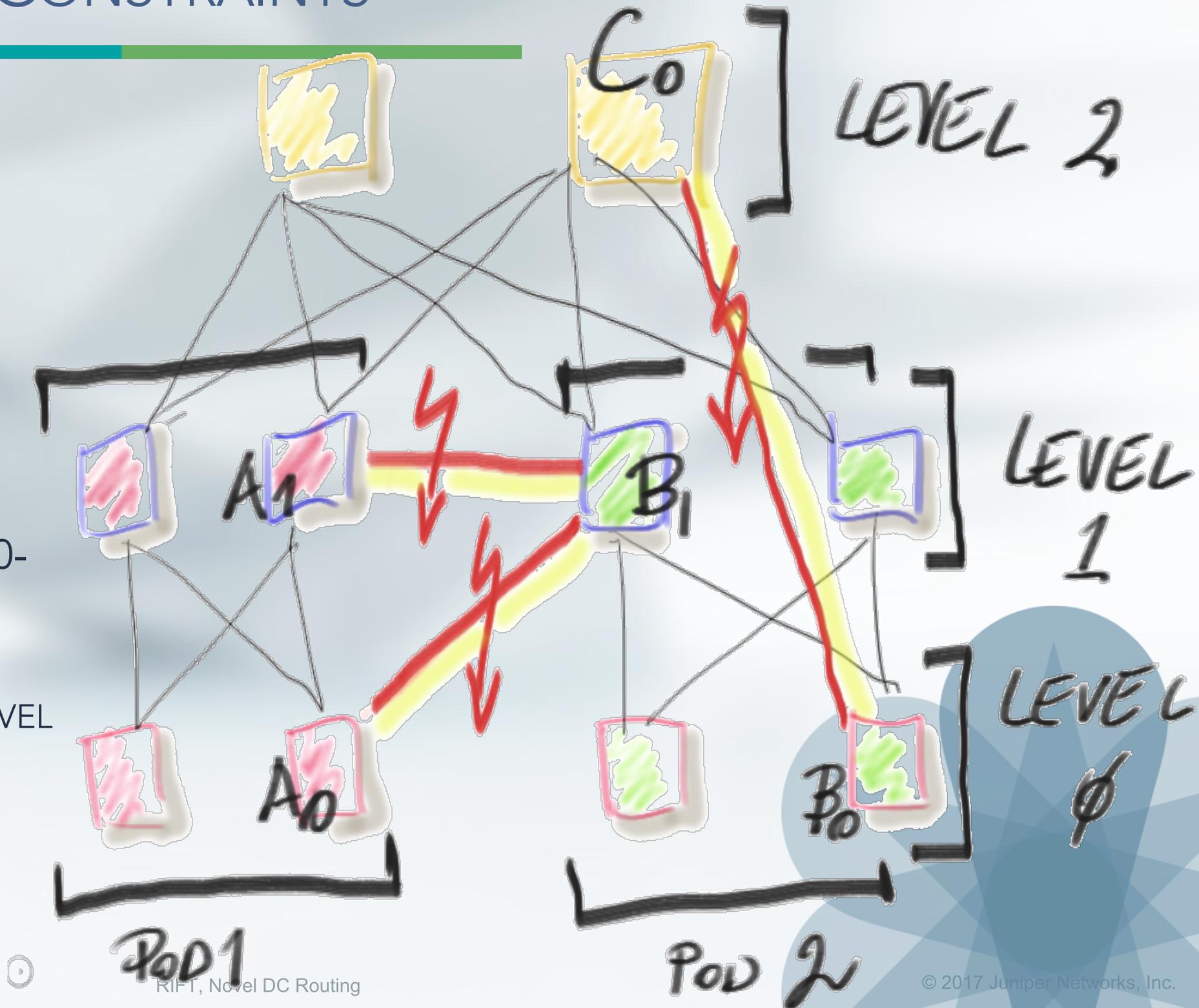
Problem / Attempted Solution	BGP modified for DC (all kind of "mods")	ISIS modified for DC (RFC7356 + "mods")	RIFT Native DC
Link Discovery/Automatic Forming of Trees/Preventing Cabling Violations	⚠	⚠	✓
Minimal Amount of Routes/Information on ToRs	⚠	⚠	✓
High Degree of ECMP (BGP needs lots knobs, memory, own-AS-path violations) and ideally NEC and LFA	⚠	✓	✓
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	✗	✓	✓
Simplest Initial Implementation	✓	✗	✗

RIFT: NOVEL DYNAMIC ROUTING ALGORITHM FOR CLOS UNDERLAY

- GENERAL CONCEPT
- AUTOMATIC CABLING CONSTRAINTS
- AUTOMATIC DISAGGREGATION ON FAILURES
- AUTOMATIC FLOODING REDUCTION
- MORE GOODIES

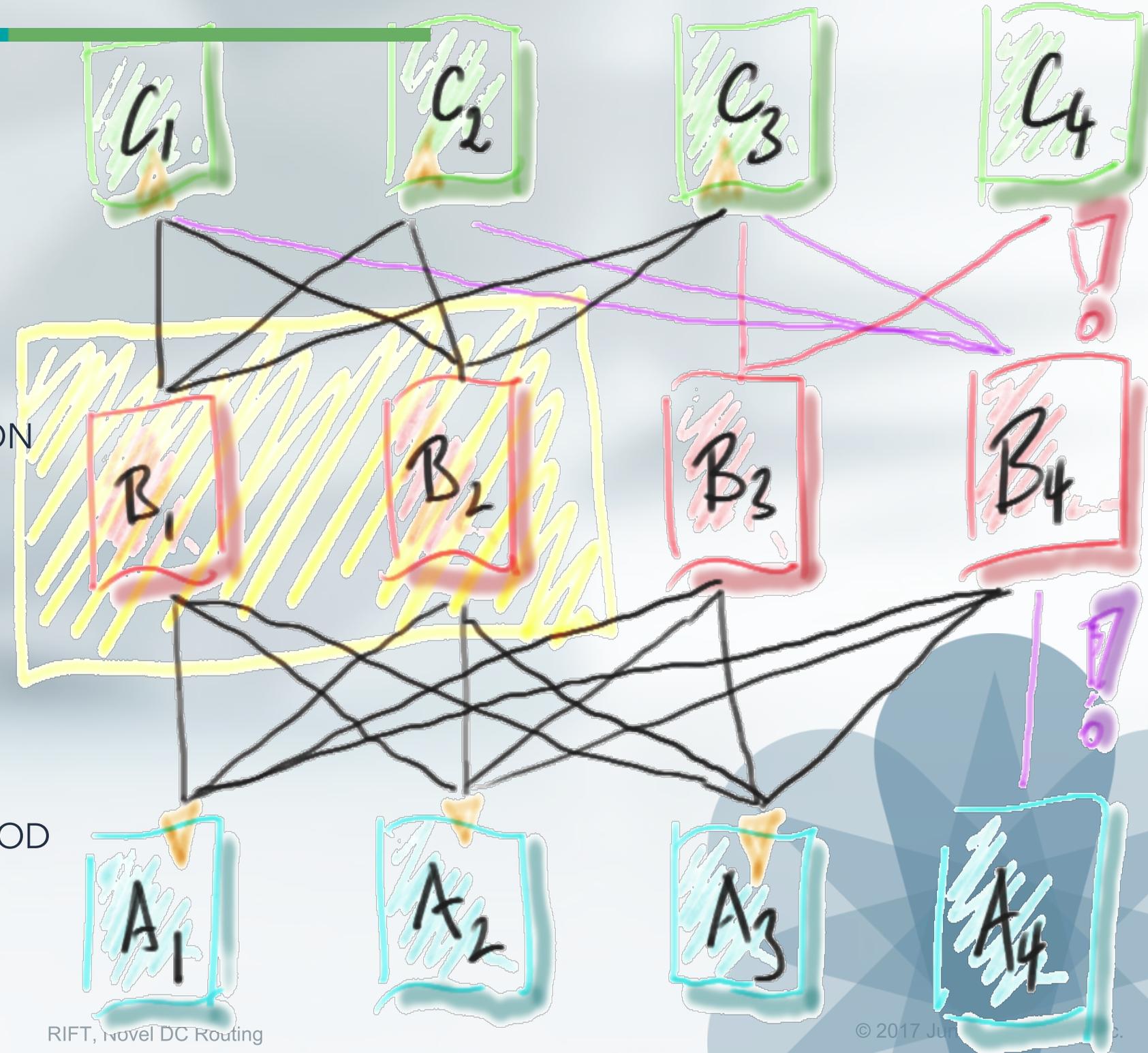
AUTOMATIC TOPOLOGY CONSTRAINTS

- LEVEL 0 = LEAF
- POD 0 = ANY POD
- AUTOMATIC REJECTION OF ADJACENCIES BASED ON MINIMUM CONFIGURATION
- A1 TO B1 FORBIDDEN DUE TO POD MISMATCH
- A0 TO B1 FORBIDDEN DUE TO POD MISMATCH (A0 ALREADY FORMED A0-A1 EVEN IF POD NOT CONFIGURED ON A0)
- B0 TO C0 FORBIDDEN BASED ON LEVEL MISMATCH
- COULD FORM OTHER TOPOLOGY VARIATIONS AS WELL



AUTOMATIC FLOODING REDUCTION

- EACH "B" NODE COMPUTES FROM REFLECTED SOUTH REPRESENTATION OF OTHER "B" NODES
 - SET OF SOUTH NEIGHBORS
 - SET OF NORTH NEIGHBORS
- NODES HAVING BOTH SETS MATCHING CONSIDER THEMSELVES "FLOOD REDUCTION GROUP" AND LOAD-BALANCE FLOODING
- FULLY DISTRIBUTED, UNSYNCHRONIZED ELECTION
- IN THIS EXAMPLE CASE B1 & B2
- EACH NODE CHOOSES BASED ON HASH COMPUTATION WHICH OTHER NODES' INFORMATION IT FORWARDS ON *FIRST* FLOOD ATTEMPT
- SIMILAR TO DF ELECTION IN EVPN BUT MUCH FASTER



POLICY GUIDED PREFIXES (PGP)

- SOUTH AND NORTH VARIANT SINCE THE "PROPAGATION DIRECTION" IS FIXED
 - AVOIDS THE "COLLIDING DIFFUSED COMPUTATION FRONTS" PROBLEMS
- PROPAGATE LIKE DISTANCE VECTOR BUT BASED ON FLOODING
 - NO NECESSITY TO BUILD SPECIALIZED UPDATES "PER PEER"
- INGRESS POLICIES CAN BE APPLIED ON PGPs
 - NO NEED FOR "REFRESHES" ON POLICY CHANGES
- USES
 - TRAFFIC ENGINEERING LIKE SR

MOREOVER

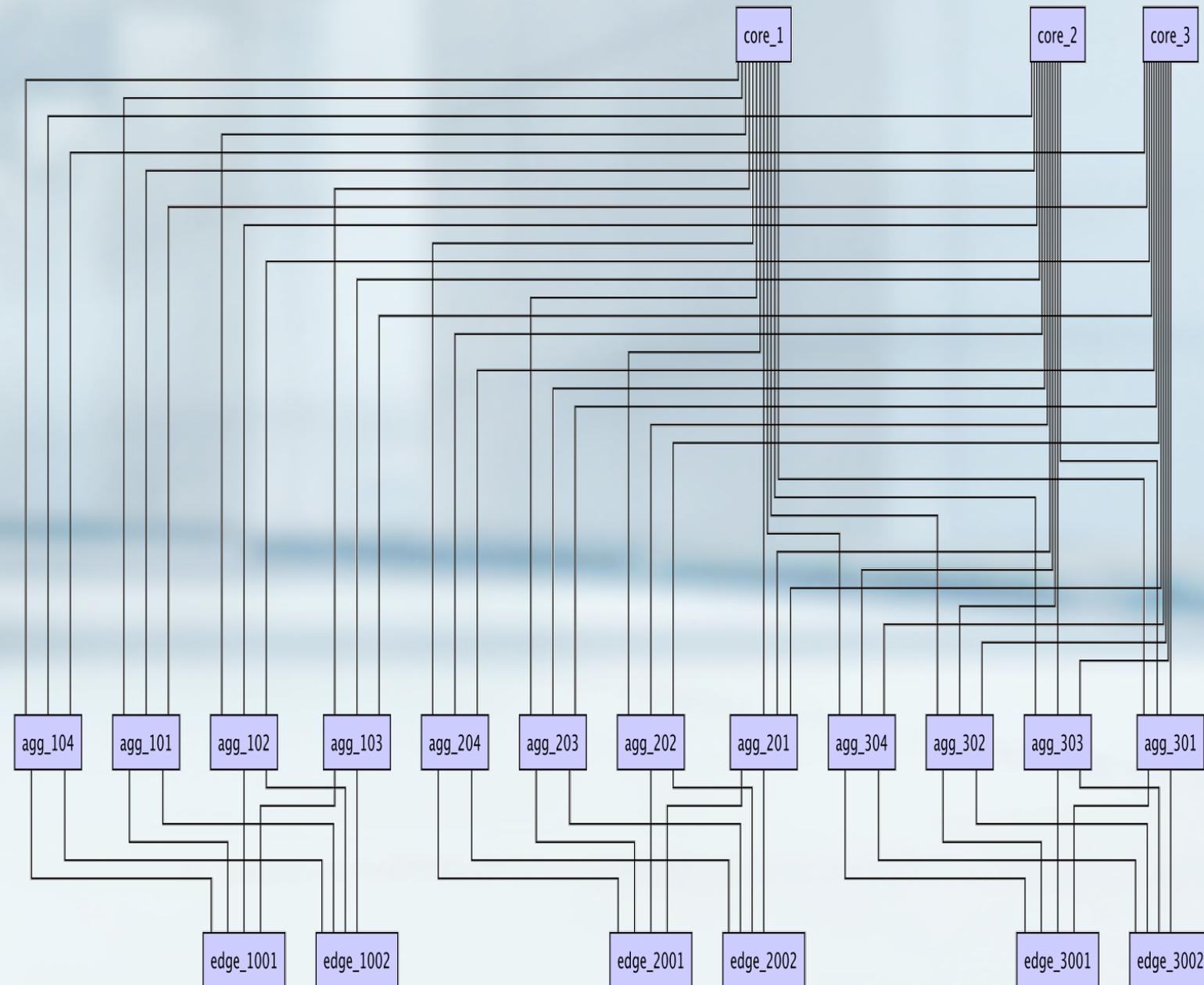
- TRAFFIC ENGINEERING, SR IS INCLUDED VIA PGP
- PACKET FORMATS ARE COMPLETELY MODEL BASED
- CHANNEL AGNOSTIC DELIVERY, COULD BE QUICK, TCP, UDP, UDT
- PREFIXES ARE MAPPED TO FLOODING ELEMENT BASED ON LOCAL HASH FUNCTIONS
 - ONE EXTREME POINT IS A PREFIX PER FLOODED ELEMENT = BGP UPDATE
- PURGING (GIVEN COMPLEXITY) IS OMITTED
- KEY-VALUE STORE IS SUPPORTED (E.G. SERVICE CONFIGURATION DURING FLOODING) INCLUDING POLICIES AND "BEST COPY TIE-BREAKING"

SUMMARY OF RIFT ADVANTAGES

- ADVANTAGES OF LINK-STATE AND DISTANCE VECTOR
 - FASTEST POSSIBLE CONVERGENCE
 - AUTOMATIC DETECTION OF TOPOLOGY
 - MINIMAL ROUTES ON TORs
 - EASY TO ACHIEVE HIGH DEGREE OF ECMP/N-ECMP
 - MINIMAL BLAST RADIUS ON FAILURES
 - FAST DE-COMMISSIONING OF NODES
- NO DISADVANTAGES OF LINK-STATE OR DISTANCE VECTOR
 - REDUCED FLOODING
 - AUTOMATIC NEIGHBOR DETECTION
- AND SOME NEITHER CAN DO
 - AUTOMATIC DISAGGREGATION ON FAILURES
 - SCOPE CONTROLLED KEY-VALUE STORE

SAMPLE COMPARISON TO IGP

- 21 NODES
- 60 LINKS
- 600 PREFIXES
- ALL RUN ON A SINGLE 4 CORES LOW END I7
- COMPARISON RIFT TO EQUIVALENT IGP
 - AVG. NODE CPU USE: 3X BETTER
 - CONVERGENCE (RIB): 4X FASTER
 - FLOODING: 4X LESS TRANSMISSIONS



THANK YOU ...

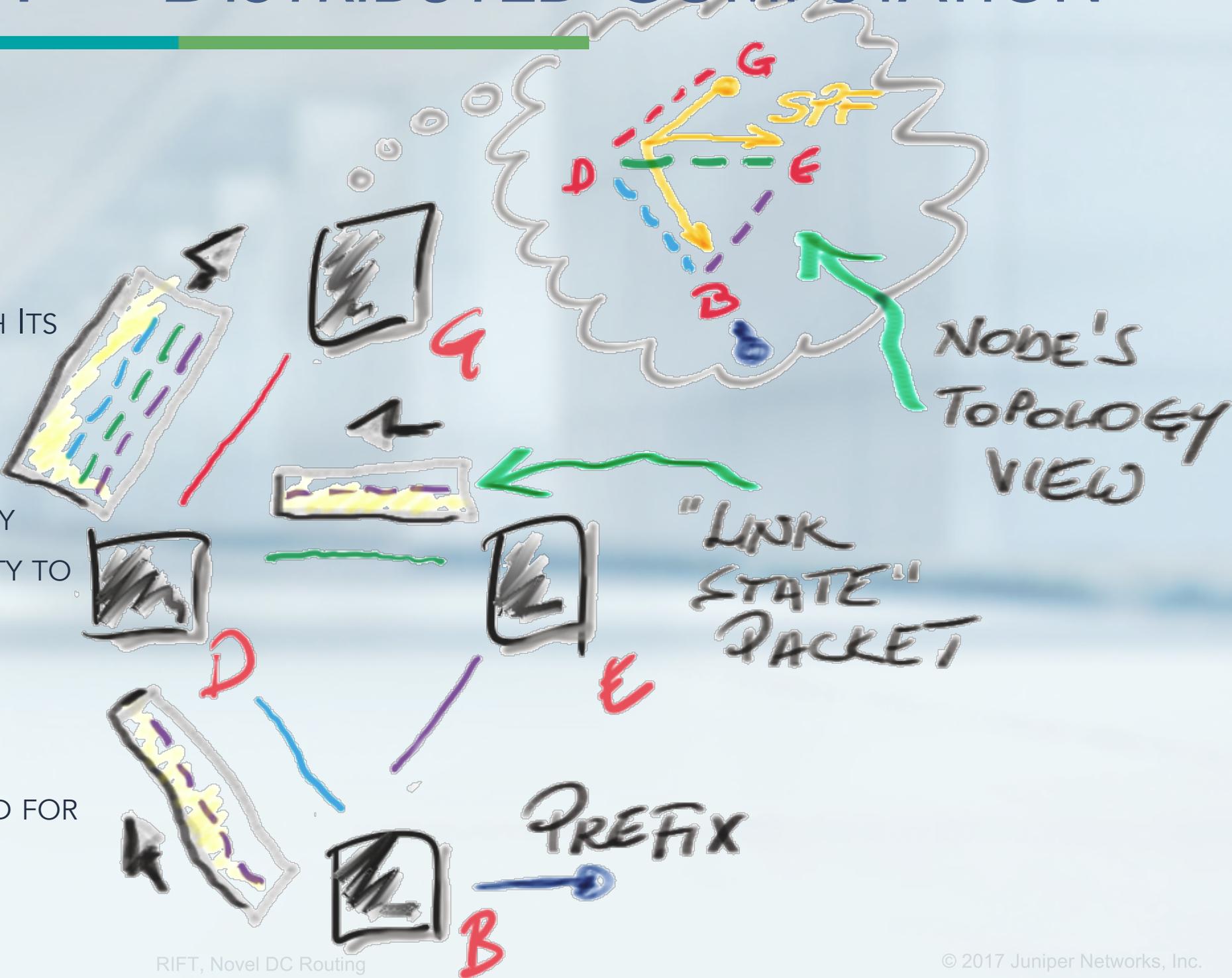
Backup Material

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
- FLOODING GENERATES EXCESSIVE LOAD FOR LARGE AVERAGE CONNECTIVITY
- PERIODIC REFRESHES



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

