

RIFT: STRAWMAN PROPOSAL OF A NOVEL DC FABRIC ROUTING PROTOCOL KEY CUSTOMERS' DECK

DRAFT-PRZYGIENDA-RIFT @ IETF

DISCLAIMERS AND EXPECTATIONS

- THIS IS A "WORKING STRAW-MAN PROPOSAL" WITH SERIOUS AMOUNT OF HIGH-PERFORMANCE PRE-PRODUCTION CODE IN PLACE
- STANDARD BOILERPLATE: NONE OF THOSE THINGS CONSTITUTE COMMITMENTS TO PRODUCT SPECIFICATIONS, OFFERINGS OR RELEASE DATES BY JUNIPER AT THIS POINT IN TIME

AGENDA

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

BLITZ OVERVIEW OF TODAY'S ROUTING

- LINK STATE & SPF
- DISTANCE/PATH VECTOR

LINK STATE AND SPF = DISTRIBUTED COMPUTATION

ୃତ

- 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)

RIFT 2017, Juniper Confidential

NODE'S TOPOLOGY VIEW

PACKET

PREFIX

DISTANCE/PATH VECTOR = DIFFUSED COMPUTATION

SIATA

REFIX

RIFT 2017, Juniper Confidential

Source

- 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^U
 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 IS VERY NOVEL AND UNPROVEN
 - 1/2 THROUGHPUT OF CLOS AT SAME COST DUE TO LOW ECMP
 - RIFT WILL 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

KOGE





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)
- YET OTHERS RUN BGP OVER IGP (TRADITIONAL ROUTING ARCHITECTURE)
- LESS THAN MORE SUCCESSFUL ATTEMPTS @ PREFIX SUMMARIZATION, MICRO- AND BLACK-HOLING

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
Peer Discovery/Automatic Forming of Trees/Preventing Cabling Violations			√
Minimal Amount of Routes/Information on ToRs	×	×	\checkmark
High Degree of ECMP (BGP needs lots knobs, memory, own-AS- path violations) and ideally NEC and LFA		\checkmark	\checkmark
Traffic Engineering by Next-Hops, Prefix Modifications	\checkmark	×	\checkmark
See All Links in Topology to Support PCE/SR		\checkmark	\checkmark
Carry Opaque Configuration Data (Key-Value) Efficiently	×		\checkmark
Take a Node out of Production Quickly and Without Disruption	×	\checkmark	√
Automatic Disaggregation on Failures to Prevent Black-Holing and Back-Hauling	×	×	\checkmark
Minimal Blast Radius on Failures (On Failure Smallest Possible Part of the Network "Shakes")	×	×	\checkmark
Fastest Possible Convergence on Failures	×	\checkmark	\checkmark
Simplest Initial Implementation	\checkmark	×	X
		DIET	2017 Junipor Con

ential

RIFT: NOVEL ROUTING ALGORITHM FOR CLOS UNDERLAY

- GENERAL CONCEPT
- AUTOMATIC TOPOLOGY CONSTRAINTS
- AUTOMATIC DISAGGREGATION
- AUTOMATIC FLOODING REDUCTION
- 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 TOPOLOGY CONSTRAINTS

LEVEL 2

¥0

POD J

LEVEL

LEVEL

RIFT 2017, Juniper Confidential

• LEVEL 0 = LEAF

14

- 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
- PROTOCOL WILL WORK AS WELL IF
 LEVEL 0 IS ALLOWED TO CONNECT TO L
 LEVEL 2 BUT OPTIMAL ROUTING WOULD
 NEED PGPS

Pool

AUTOMATIC DE-AGGREGATION



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 Matching Sets Consider Themselves "Flood Reduction Group"
- FULLY DISTRIBUTED, UNSYNCHRONIZED ELECTION
- IN OUR 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



MOREOVER

- TRAFFIC ENGINEERING VIA "FLOODED DV OVERLAY" WITH POLICIES
- COMPLETELY MODEL BASED PACKET FORMATS
- CHANNEL AGNOSTIC DELIVERY, COULD BE QUICK, TCP, UDP, UDT
- PREFIXES TO TOPOLOGY ELEMENT MAPPING BASED ON HASH FUNCTIONS LOCAL TO EACH NODE
 - ONE EXTREME POINT IS PREFIX PER FLOODED ELEMENT = BGP UPDATE
- PURGING (GIVEN COMPLEXITY) IS OMITTED
- POLICY CONTROLLED KEY-VALUE STORE SUPPORT

SUMMARY OF RIFT ADVANTAGES

- ADVANTAGES OF LINK-STATE AND DISTANCE VECTOR
 - FASTEST POSSIBLE CONVERGENCE
 - AUTOMATIC DETECTION OF TOPOLOGY
 - MINIMAL ROUTES ON TORS
 - HIGH DEGREE OF ECMP
 - MINIMAL BLAST RADIUS ON FAILURES
 - FAST DE-COMISSIONING OF NODES
 - MAXIMUM PROPAGATION SPEED WITH FLEXIBLE # PREFIXES IN AN UPDATE

- NO DISADVANTAGES OF LINK-STATE OR DISTANCE VECTOR
 - REDUCED FLOODING
 - AUTOMATIC NEIGHBOR DETECTION
- AND SOME NEITHER CAN DO
 - AUTOMATIC DISAGGREGATION ON FAILURES
 - KEY-VALUE STORE

EARLY JUNIPER IMPLEMENTATION

- PRE-PRODUCTION CODE
- FLOODING/HELLO/SPF/RIB/FLOOD REDUCTION IMPLEMENTED
- OPTIONAL REDIS PERSISTENCY FOR LINK-STATE DB AND STATISTICS IMPLEMENTED
- VERY HEAVILY MULTI-THREADED TO FULLY UTILIZE MODERN CPU ARCHITECTURES
 - MEMORY AND THREAD SAFE UNLIKE C
- OWN TOPOLOGY ENVIRONMENT FOR HEAVY-DUTY TESTING
- EASILY PORTABLE ACROSS ALL MODERN OS PLATFORMS

INDICATIVE NUMBERS FOR RIFT: TOPOLOGY



RIFT 2017, Juniper Confidential

20

INDICATIVE NUMBERS FOR RIFT IMPLEMENTATION

٠

- MACPRO (LOW POWER I7 WITH 4 REAL CORES)
- 21 NODES/60 LINKS/600 PREFIXES EXECUTED
- ~480 THREADS
- FLOODING ON HYSTERESIS
 - AROUND 10+K TIES/SEC PER ADJACENCY

- CONVERGENCE FROM COLD START
 - TOTAL CPU USE (4 CORES I7)
 - FLAT: 11SECS ~ 500 MSECS/NODE
 - RIFT: 3.2SECS ~ 150 MSECS/NODE
 - RIB CONVERGENCE FOR CORE SUPER-SPINES
 - FLAT: 900MSEC
 - RIFT: ~250MSEC
 - AMOUNT OF FLOODING (NO FLOOD REDUCTION APPLIED YET)
 - FLAT: ~13K TIES, 96K TOTAL TRANSMISSIONS
 - RIFT: ~5K TIES, 21K TOTAL TRANSMISSIONS

SINGLE LINK FLAP (CORE 1 TO AGG 104)

- FLAT: AVG: ~60 MSECS, MAX: ~75 MSECS
- RIFT: AVG: ~35 MSECS, MAX: ~70 MSECS

ONGOING/FUTURE WORK

- AUTOMATIC DE-AGGREGATION
- NORTHBOUND COMPUTATION
- LARGE SCALE RUNS ON SERVER FARMS
- PREFIX GUIDED PREFIXES
- INTEGRATION WITH ONE OF OUR PLATFORMS
- TI[DR]ES ON TOS (OR DIFFERENT PORT)
- KEY CUSTOMER ENGAGEMENTS
 - REQUIREMENTS/USE CASES
 - DRAFT CO-AUTHORSHIP DESIRED



THANKS