

# RIFT: TOWARDS SELF-DRIVING IP FABRICS WITH ZERO OPEX ROUTING

---

TONY PRZYGIENDA PRZ@JUNIPER.NET

DRAFT-PRZYGIENDA-RIFT @ IETF

# DISCLAIMERS AND EXPECTATIONS

---

NONE OF THOSE THINGS CONSTITUTE COMMITMENTS TO PRODUCT SPECIFICATIONS, OFFERINGS OR RELEASE DATES BY JUNIPER AT THIS POINT IN TIME

# AGENDA

---

- EVERYTHING IN A NUTSHELL
- FABRIC TOPOLOGY
- BLITZ OVERVIEW OF TODAY'S ROUTING
- "IP FABRIC ROUTING" IS A SPECIALIZED PROBLEM
- RIFT: ZERO OPEX ROUTING FOR SELF-DRIVING IP CLOS FABRICS

# IN A NUTSHELL

---

- DATA CENTERS ARE A STRATEGIC ASSET FOR LARGE CORPORATIONS AND OPERATORS
  - IF THEY DON'T OWN THEIR CRITICAL DATA PROCESSING, THEIR DATA OWNER MAY OWN THEM IN THE FUTURE
- BUILDING DATA CENTERS NECESSITATES BUILDING FABRICS
  - LARGE FABRICS ARE BECOMING IP ONLY AND NEED TO BE ROUTED
- ROUTING OPEX FOR IP FABRICS IS SIGNIFICANT
  - TECHNICALLY COMPLICATED, NEITHER EFFICIENT NOR ROBUST
  - FEW AVAILABLE EXPERTS IN THE FIELD
  - “SELF DRIVING” FABRICS ARE NEEDED
    - ZERO CONFIGURATION, RESILIENCE, MAXIMUM GOODPUT, AS “SELF DRIVING” AS POSSIBLE
- JUNIPER IS WORKING ON OPEN STANDARD FOR ZERO OPEX, SELF-DRIVING IP FABRICS ROUTING
- ACCIDENTALLY, “EXPLODED CHASSIS BACKPLANE” PRESENTS THE SAME PROBLEM

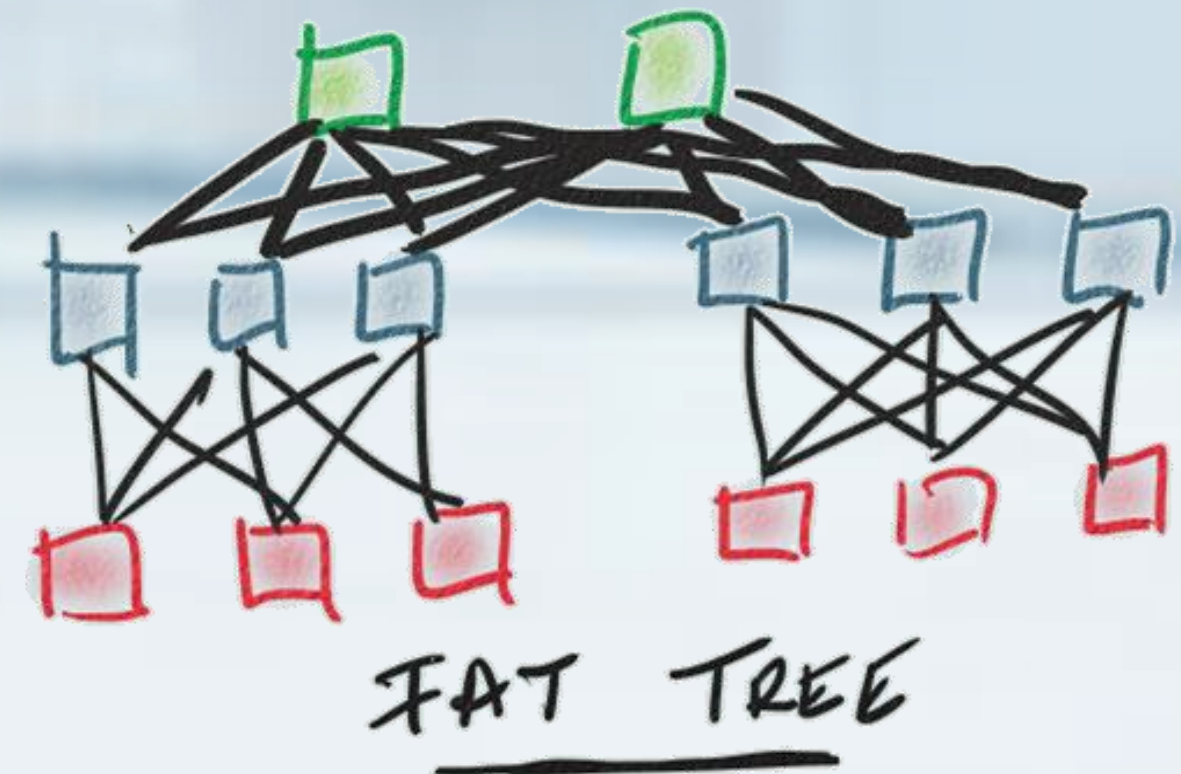
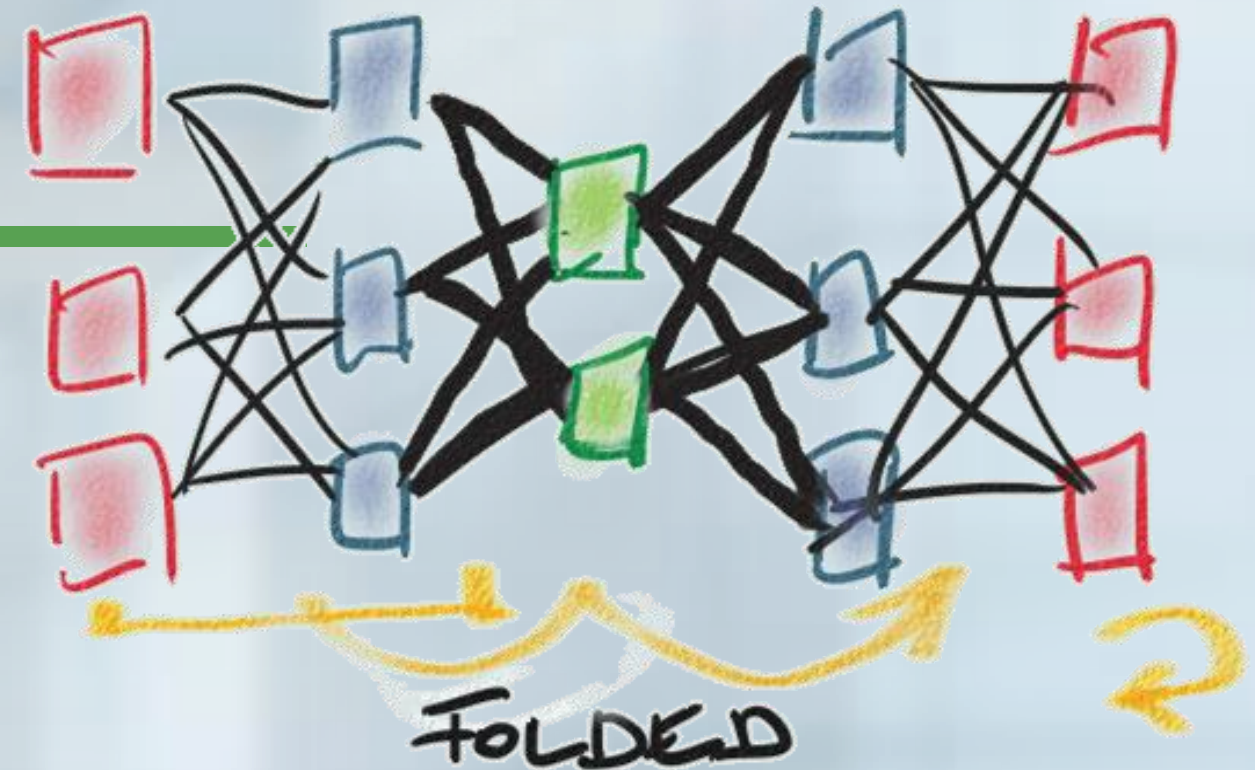
# FABRIC: A SPECIALIZED TOPOLOGY

---

- CLOS TOPOLOGIES ARE DOMINANT TODAY
- CURRENT STATE OF IP FABRIC ROUTING AFFAIRS
- REQUIREMENTS MATRIX FOR SELF-DRIVING IP FABRIC ROUTING

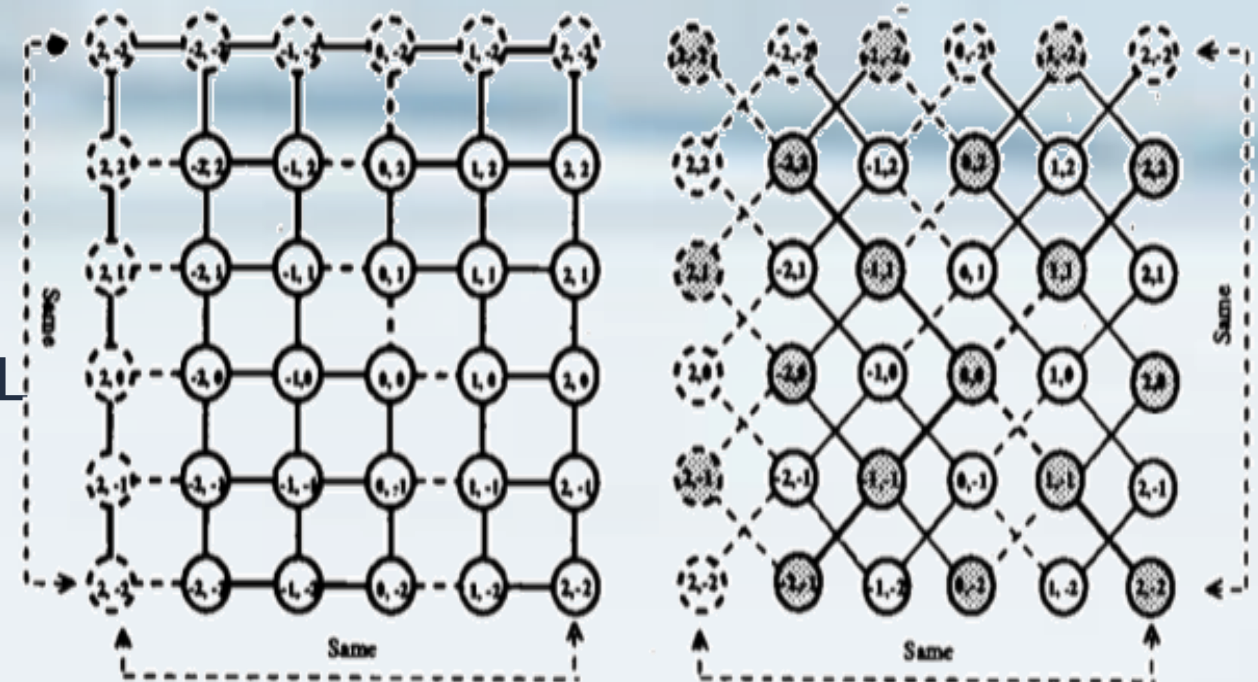
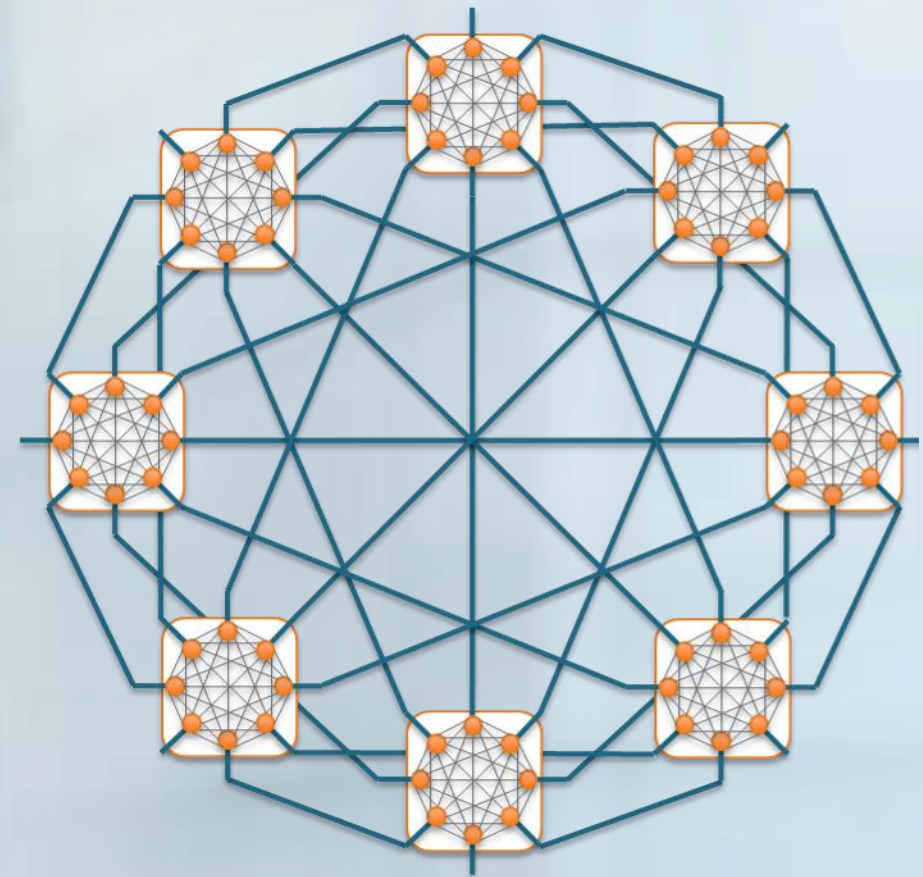
# CLOS TOPOLOGIES

- 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



# WHY NOT SOMETHING ELSE ?

- TOROIDAL [AND DIAGONAL] MESHES HAVE LONG PATHS, SMALL BISECTION WIDTH AND POOR BLOCKING PROPERTIES
- DRAGONFLY IS VERY NOVEL AND UNPROVEN
  - SEEMINGLY  $\frac{1}{2}$  THROUGHPUT OF CLOS AT SAME CAPACITY DUE TO LOW ECMP
  - OUR SUGGESTION SHOULD WORK WELL IN A PRACTICAL MODIFICATION (ONE LEVEL CLOS AND DRAGONFLY CORE) IF NECESSARY



# BLITZ OVERVIEW OF TRADITIONAL ROUTING

---

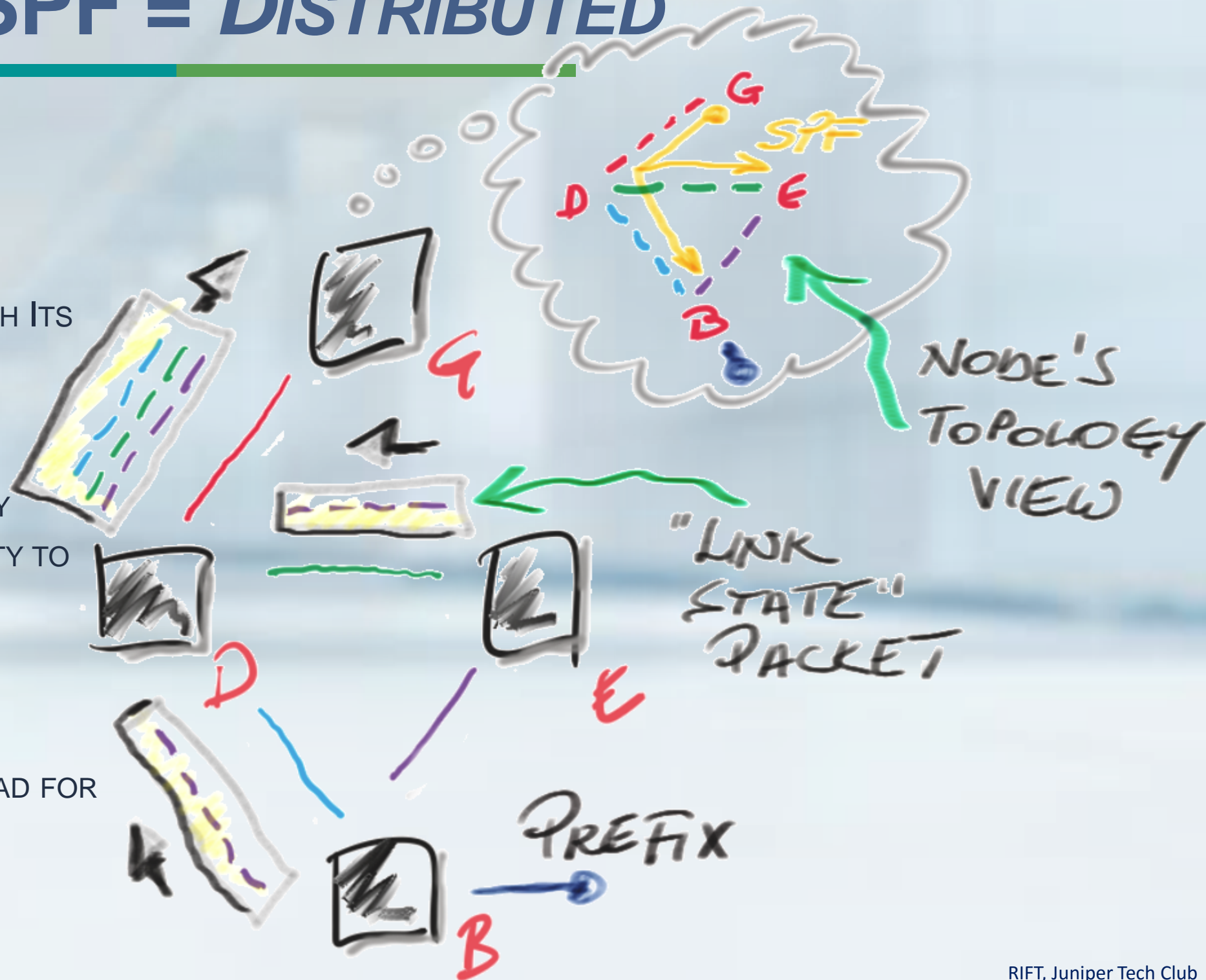
- LINK STATE & SHORTEST PATH FIRST
- DISTANCE & PATH VECTOR



# LINK STATE AND SPF = *DISTRIBUTED*

## COMPUTATION

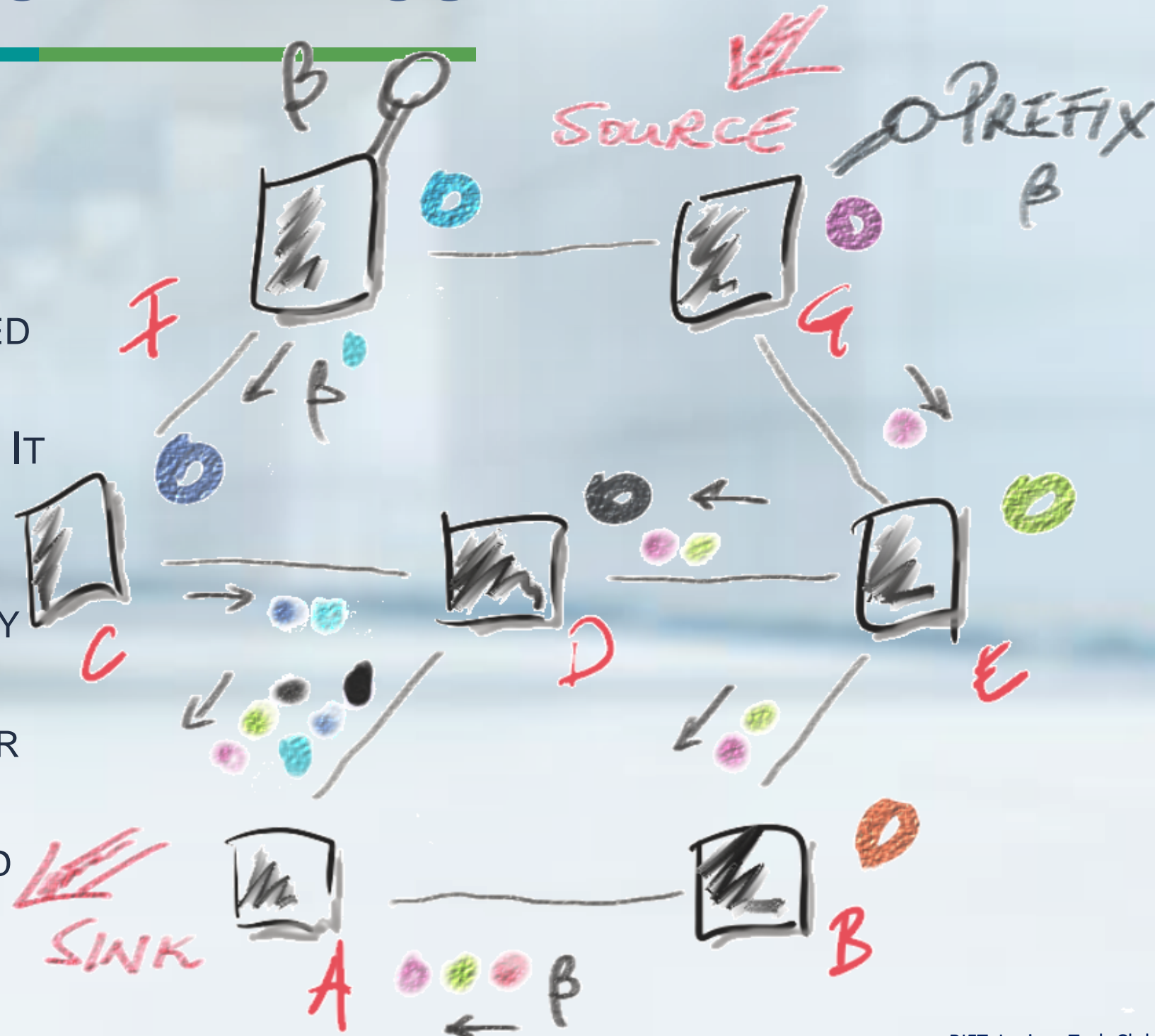
- TOPOLOGY ELEMENTS
  - NODES
  - LINKS
  - PREFIXES
- EACH NODE ORIGINATES PACKETS WITH ITS TOPOLOGY ELEMENTS
- PACKETS ARE "FLOODED"
- "NEWEST" VERSION WINS
- EACH NODE "SEES" WHOLE TOPOLOGY
- EACH NODE "COMPUTES" REACHABILITY TO EVERYWHERE
- CONVERSION IS VERY FAST
- EVERY LINK FAILURE SHAKES WHOLE NETWORK (MODULO AREAS)
- FLOODING GENERATES EXCESSIVE LOAD FOR LARGE ADJACENCY FAN-OUTS
- 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 ENFORCING “POLICY” RATHER THAN PROVIDE “MAX REACHABILITY”
- SCALES WHEN PROPERLY IMPLEMENTED TO MUCH HIGHER # OF ROUTES THAN LINK-STATE



# CURRENT STATE OF ROUTING IN DC FABRICS

- 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
- MANY LARGE CORPORATIONS RUN “FLAT” IGP (ISIS OR OSPF)
- YET OTHERS RUN BGP OVER IGP (TRADITIONAL ROUTING ARCHITECTURE)
- LESS THAN MORE SUCCESSFUL ATTEMPTS @ PREFIX SUMMARIZATION, CONTROL OF MICRO- AND BLACK-HOLING

# REQUIREMENTS BREAKDOWN (RFC7938+) FOR A “ZERO OPEX FABRIC”

Problem / Attempted Solution	BGP modified for DC	ISIS modified for DC	RIFT
01. As Close to Zero Necessary Configuration as Possible (Contradicts 02)	X	X	✓(*)
02. Peer Discovery/Automatic Forming of Trees/Preventing Cabling Violations (Contradicts 01)	⚠️	⚠️	✓
03. Minimal Amount of Routes/Information on ToRs	X	X	✓
04. High Degree of ECMP (BGP needs lots knobs, memory, own-AS-path violations) and ideally NEC and LFA	⚠️	✓	✓
05. Traffic Engineering by Next-Hops, Prefix Modifications	✓	X	✓
06. See All Links in Topology to Support PCE/SR	⚠️	✓	✓
07. Carry Opaque Configuration Data (Key-Value) Efficiently	X	⚠️	✓
08. Take a Node out of Production Quickly and Without Disruption	X	✓	✓
09. Automatic Disaggregation on Failures to Prevent Black-Holing and Back-Hauling	X	X	✓
10. Minimal Blast Radius on Failures (On Failure Smallest Possible Part of the Network “Shakes”)	X	X	✓
11. Fastest Possible Convergence on Failures	X	✓	✓

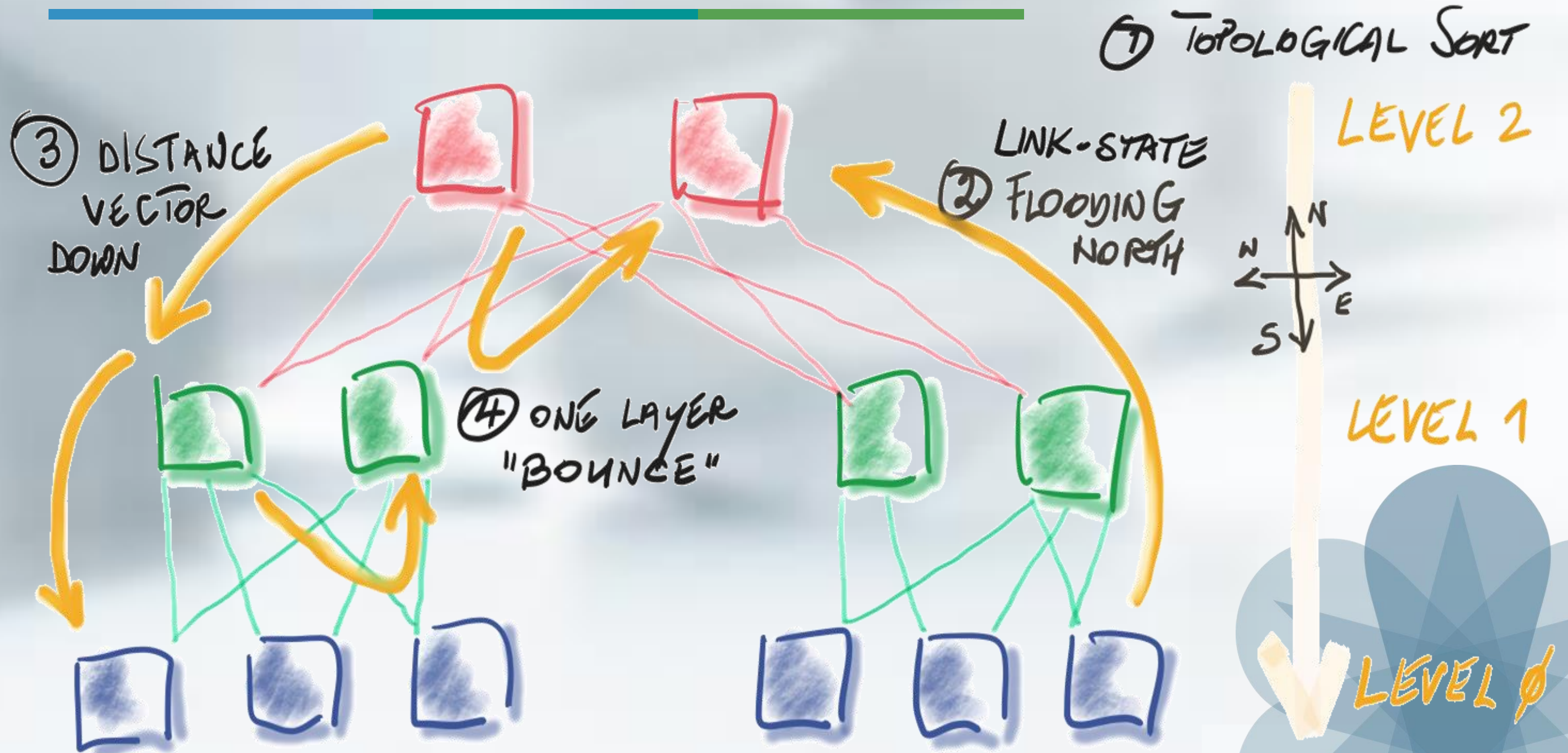
# RIFT: SELF DRIVING ROUTING ALGORITHM FOR CLOS UNDERLAY

- GENERAL CONCEPT
- AUTOMATIC MIS-CABLING CONSTRAINTS
- AUTOMATIC DISAGGREGATION
- OPTIONAL HORIZONTAL LINKS
- AND MORE BEYOND THAT

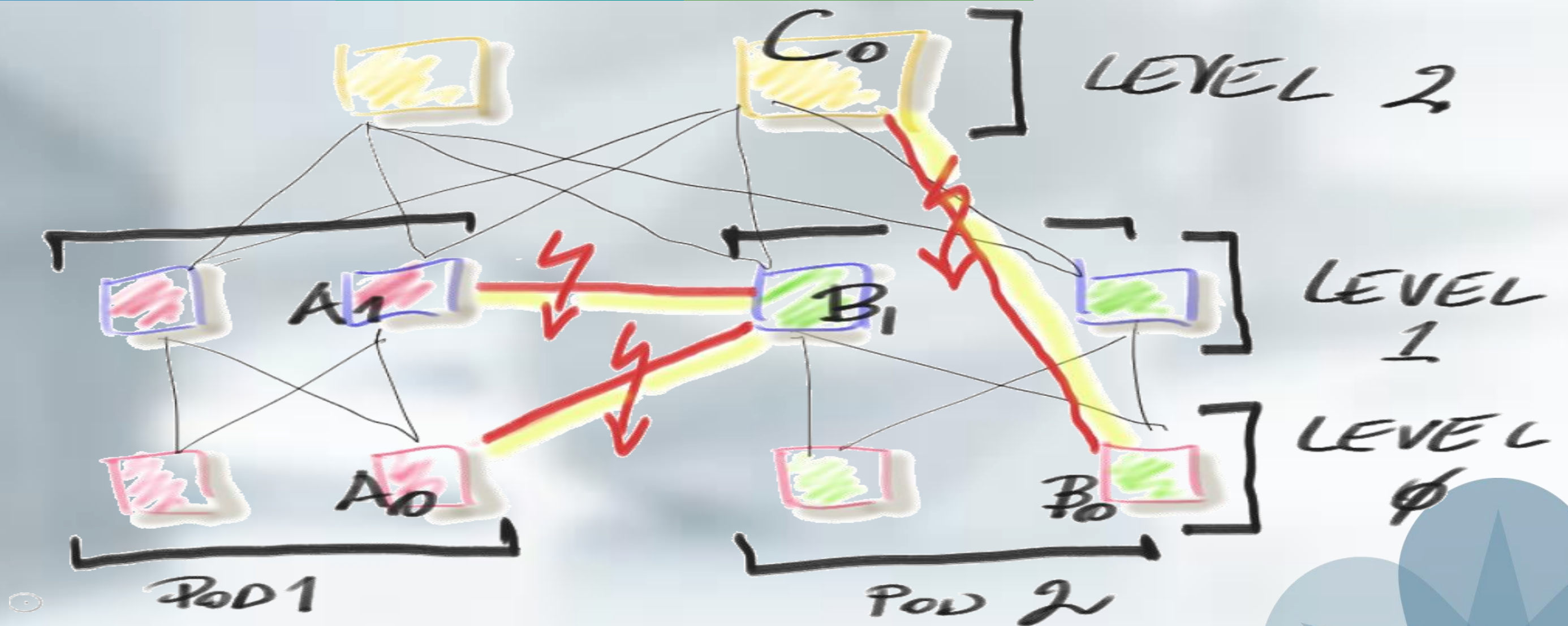
*“Just because the standard provides a cliff in front of you, you are not necessarily required to jump off it.”*

*— Norman Diamond*

# LINK-STATE UP, DISTANCE VECTOR DOWN & BOUNCE

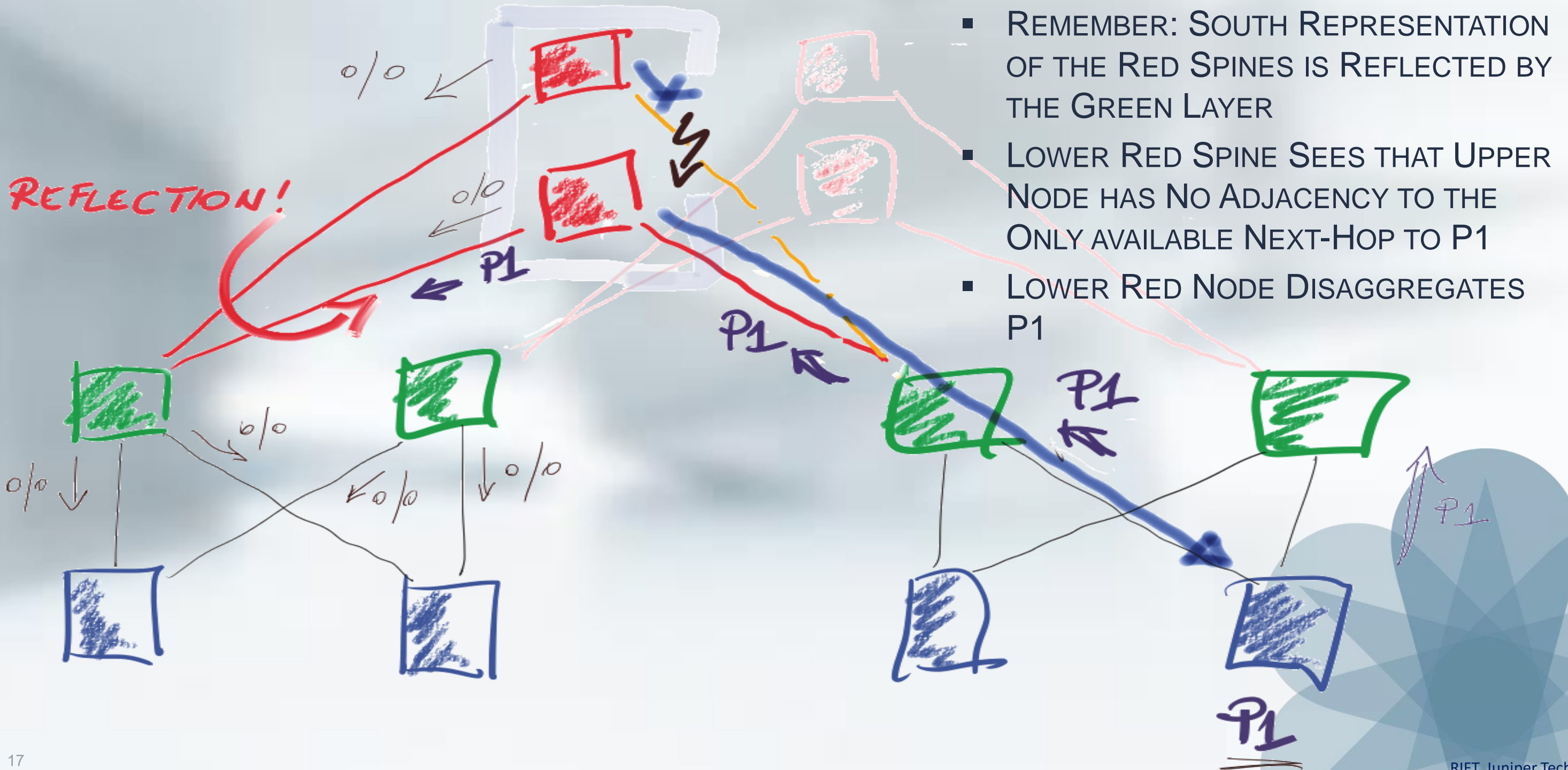


# AUTOMATIC MIS-CABLING CONSTRAINTS



- AUTOMATIC REJECTION OF ADJACENCIES BASED ON MINIMUM CONFIGURATION
- PROTOCOL WILL WORK AS WELL IF LEVEL 0 IS ALLOWED TO CONNECT TO LEVEL 2 BUT OPTIMAL ROUTING WOULD NEED LARGER FIBS ON LEAFS

# AUTOMATIC DISAGGREGATION

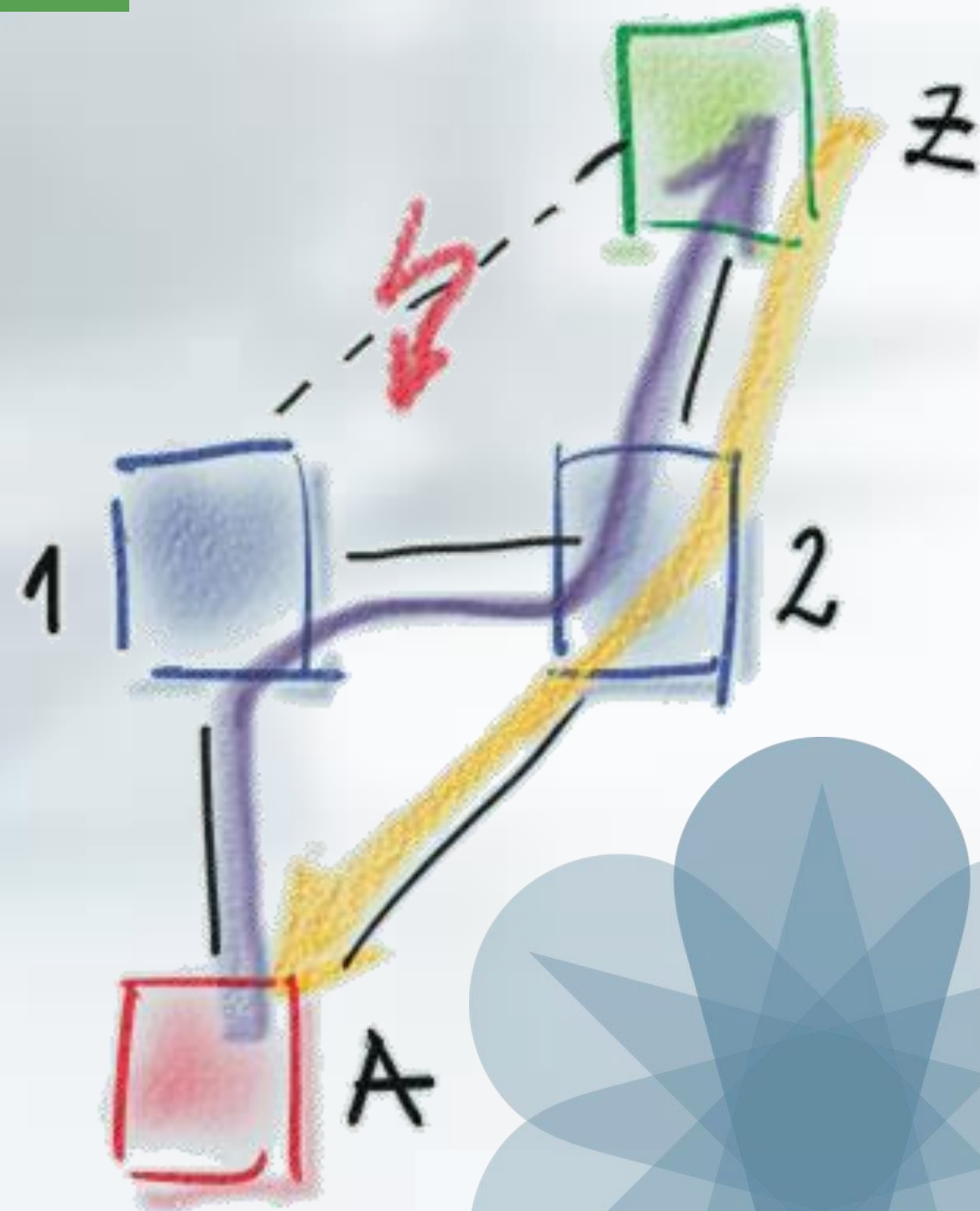




# OPTIONAL HORIZONTAL LINKS FOR LAYER

## PROTECTION

- LEVELS CAN INSTALL OPTIONAL HORIZONTAL LINKS
- LEVEL 0 IS SPECIAL:
  - LEAF-2-LEAF CONNECTION THAT CANNOT BE USED EXCEPT FOR LEAF-2-LEAF TRAFFIC
- LEVEL > 0 USES HORIZONTAL LINKS FOR FAILURE PROTECTION ONLY
  - SINGLE NODE PROTECTION: NODE THAT LOST NORTHBOUND LINKS BUT HAS NEIGHBORS THAT CAN REACH HIGHER LAYERS USES THE HORIZONTAL LINK
  - N:N-1 PROTECTION: FULL MESH IN A LEVEL CAN PROVIDE UP TO N-2 NORTHBOUND PROTECTION



# RIFT Does On Top

---

- AUTOMATIC FLOOD REDUCTION
- LEAF-TO-LEAF BI-DIRECTIONAL SHORTCUTS
- POSSIBLE TRAFFIC ENGINEERING VIA “FLOODED DV OVERLAY” WITH POLICIES
- COMPLETELY MODEL BASED PACKET FORMATS
- CHANNEL AGNOSTIC DELIVERY, COULD BE QUICK, TCP, UDP
- 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

---

- OPEN IETF STANDARD
- **ADVANTAGES OF BOTH 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
- **NO DISADVANTAGES OF NEITHER LINK-STATE NOR DISTANCE VECTOR**
  - REDUCED FLOODING
  - AUTOMATIC NEIGHBOR DETECTION
- **UNIQUE RIFT ADVANTAGES**
  - AUTOMATIC DISAGGREGATION ON FAILURES
  - KEY-VALUE STORE
  - HORIZONTAL LINKS USED FOR PROTECTION
  - MINIMAL BLAST RADIUS ON FAILURES
  - CAN UTILIZE **ALL** PATHS THROUGH FABRIC WITHOUT LOOPING

# IS THERE MORE THAN BITS OVER POWERPOINT



- DETAILED DRAFT IN IETF
  - [HTTPS://DATATRACKER.IETF.ORG/DOC/DRAFT-PRZYGIENDA-RIFT/](https://datatracker.ietf.org/doc/draft-przygienda-rift/)
- PRE-PRODUCTION CODE AVAILABLE UNDER NDA
  - PLEASE TALK TO YOUR FRIENDLY SYSTEMS OR RESIDENT ENGINEER
- SINGAPORE IETF IS HOSTING A “DC ROUTING BoF” SESSION

# THANKS

