

ROB SHAKIR END-TO-END NETWORK ARCHITECT CHIEF NETWORK ARCHITECT'S OFFICE, BT.

PARIS, FRANCE – MARCH 2014.

HI!



ROB SHAKIR END-TO-END NETWORK ARCHITECT BT plc



TECHNOLOGY STRATEGY: LEGACY TO STRATEGIC NETWORK EVOLUTION: INTRODUCE NEW SERVICES & CAPABILITIES.

LTE Mobile

IN THE BEGINNING.

catalyst2

Small peering network offering content hosting and managed services. 3 UK PoPs, 3 staff, <50Mbps traffic!

DISCOVERING LEGACY.



NSP offering Internet, data, voice and hosting services – lots of acquisitions and legacy!

20+ EU PoPs, 250 staff, 10s of Gbps traffic.

GOING GLOBAL – TELCO #1.



Global telco offering transport, IP/ MPLS, voice, services – focus on multi-service & IP convergence. 200+ PoPs, ~10,000 staff, 100s of Gbps traffic.

CURRENTLY – GLOBAL TELCO

Many strategic and legacy networks with O(10,000s) elements – spread over ~6,000 PoPs in 180 countries.

Delivering private and public Internet services – edge capacity of terabits/second peak.



COMMON THEME: EVERYTHING IS IP.





APPLICATIONS EXPECT LEGACY NETWORKING BEHAVIOUR – INTRODUCING NEW PATH ROUTING AND PERFORMANCE REQUIREMENTS TO IP.

COMMON THEME: MULTI-SERVICE NETWORKS.



MULTI-PROTOCOL LABEL SWITCHING? Label imposed to indicate







FORWARDING BASED ON LABEL INFORMATION RATHER THAN IP DESTINATION – IP TUNNELLED INSIDE MPLS PATH (FEC). LABEL USED TO INDICATE PACKET'S CONTEXT – ALLOWS MULTIPLE SERVICES TO BE USED – BOTH L3 + L2.

BASIC MPLS FORWARDING.



TODAY'S BASIC IP/MPLS NETWORK ANATOMY.



BASIC NETWORK – HOW DO WE DISCOVER WHICH LINKS ARE UP/DOWN AND WHICH DESTINATIONS ARE REACHABLE THROUGH ADJACENT NODES?

TODAY'S BASIC IP/MPLS NETWORK ANATOMY – IGP.



INTERIOR GATEWAY PROTOCOL (IGP) – TYPICALLY OSPF OR IS-IS – PROVIDES INFORMATION ABOUT REACHABILITY BETWEEN NODES, ALLOWING SHORTEST-PATH CALCULATIONS BASED ON ADVERTISED COST.

TODAY'S BASIC IP/MPLS NETWORK ANATOMY – IGP.



FORWARDING THROUGH THE NETWORK IS THEN BASED ON SHORTEST PATH INFORMATION ONLY – SINGLE SET OF METRICS FOR THE NETWORK.

TODAY'S BASIC IP/MPLS NETWORK ANATOMY – LDP.



ALLOWS MPLS FORWARDING ALONG THE IGP SHORTEST PATH – BUT INCREASES THE NUMBER OF PROTOCOLS DEPLOYED AND OPERATIONAL COMPLEXITY.

SELECTING A NON-SHORTEST PATH.



Pink path is divergent from lowest cost path (cost = 25, rather than 20)

FOR SOME MULTI-SERVICE APPLICATIONS – E.G., PATH DISJOINTNESS – WE NEED TO SELECT A PATH WHICH IS NOT THE IGP SHORTEST PATH.

SELECTING A NON-SHORTEST PATH – RSVP-TE.



FURTHER PROTOCOL INTRODUCED TO ALLOW FOR EXPLICIT PATHS – WHICH MAY BE A SMALL SUBSET OF THE TRAFFIC CARRIED ON THE NETWORK.



HEAD-END (INITIATOR) OF LABEL SWITCHED PATH (LSP) CALCULATES PATH THROUGH THE NETWORK THAT IS REQUIRED - LDP/IGP DOES NOT GUARANTEE THIS PATH.



HEAD-END DEVICE BUILDS A MESSAGE TO SIGNAL THIS LSP – INCLUDING ANY CONSTRAINTS REQUIRED FOR THE NEW TUNNEL.



HEAD-END DEVICE BUILDS A MESSAGE TO SIGNAL THIS LSP – INCLUDING ANY CONSTRAINTS REQUIRED FOR THE NEW TUNNEL.



WHEN DESTINATION RECEIVES PATH MESSAGE, IT CREATES A RESERVATION MESSAGE PROVIDING LABEL INFORMATION FOR THE LSP.



RESERVATION MESSAGE CREATES STATE ON EACH NODE ALONG THE PATH – AND REPORTS LABEL INFORMATION BACK TO THE HEAD-END FOR FORWARDING.



AS WELL AS DIFFERING SET UP PROCEDURES, THE RESERVATIONS MUST BE REFRESHED PERIODICALLY (SOFT-STATE PROTOCOL REQUIREMENT).



DURING LINK FAILURES – HEAD-END IS NOT REQUIRED TO RE-CONVERGE PATHS – BUT RATHER NODES ADJACENT TO THE FAILURE MUST REPORT PATH FAILURE – RESULTING IN RE-SIGNALLING PROCESS RE-STARTING.

APPLICATIONS OF RSVP-TE LSPS.



FAST RE-ROUTE: EXPLICIT PATH AVOIDING A LINK OR NODE TO BE PROTECTED WHICH CAN BE USED TO PROVIDE FAST RESTORATION OF A PATH.

APPLICATIONS OF RSVP-TE LSPS.



DISJOINT PATHS: RED/BLUE TOPOLOGIES CAN BE SIGNALLED BASED ON LSPS WHICH TRAVERSE PARTICULAR SETS OF LINKS.

APPLICATIONS OF RSVP-TE LSPS.



TRAFFIC ENGINEERING: PLACEMENT OF PATHS ACCORDING TO AVAILABLE RESOURCES SUCH AS BANDWIDTH OR LATENCY.

RSVP-TE IN THE CONTEXT OF RFC5218.

Looking at RSVP-TE in the context of 5218 – we can see where these solutions were intended to fit – and the possible expansions out to the generic multi-service IP/MPLS network use-cases discussed originally.





SCALING SOFT-STATE PROTOCOLS:

ALTHOUGH RSVP-TE CAN SIGNAL THESE PATHS – SCALE IS LIMITED BY THE NUMBER OF PATHS THAT CAN BE REFRESHED WITHIN THE INTERVAL (OR RATHER, IN THE SCHEDULED CPU CYCLE TIME).



RELATIVELY EASY TO SOLVE – BUT REQUIRED ADDITIONS TO PROTOCOL. REFRESH ALL LSPS WITHIN A SINGLE MESSAGE – REDUCES NUMBER OF MESSAGES THAT MUST BE GENERATED.

But the soft state problem is not wholly solved outside of steady-state operation – for example, look at a large failure on a mid-point carrying many LSPs (quite common due to sub-sea cable connectivity)!



PATHTEAR PATHTEAR PATH PATHTEAR 4 RESV 3 PATH 2 PATHTEAR RESV PATH PATH PATH PATHTEAR PATHTFAR PATHTEAR PATH PATH

RETRY TIME

1. GLOBAL REVERT PATH.

Head-end transmits a Path message for LSP Tunnel ID = M, LSP ID = N.

2. PATHTEAR SENT AFTER RETRY INTERVAL. Head-end tears down Tunnel ID = M, LSP ID = N- retry period expired.

3. HEAD-END RE-SENDS PATH. HE tries to signal new LSP – Tunnel ID = M, LSP ID = N+1.

4. MID-POINT RECEIVES RESV BACK.

Resv received for Tunnel ID = M, LSP ID = \underline{N} – out of date!

CONTINUAL LOOP!

Midpoint never processes the 'current' LSP ID. Results in at least two further messages in the queue. Results in a "Snowball" effect.

RSVP-TE had an incremental deployment advantage and solves real operator challenges...

But...comes with <u>scalability limits which can result</u> in significant fragility being introduced into real <u>networks</u> – this is definitely **negative net value**!



CHALLENGE 2: OPERATIONAL COMPLEXITY.



PATH PLACEMENT: WHERE SHOULD A PARTICULAR LSP BE PLACED WITHIN THE NETWORK TO ENSURE NO SHARED RISKS, AND SUCH THAT WE MAXIMISE ROBUSTNESS?

CHALLENGE 2: OPERATIONAL COMPLEXITY.



PATH PLACEMENT FOR TRAFFIC ENGINEERING: WHERE WAS A PARTICULAR PATH PLACED (NOT JUST RELATED TO THE AVAILABLE LINKS) AT A PARTICULAR TIME? WHICH SERVICES WERE IMPACTED BY A CERTAIN FAILURE?

CHALLENGE 2: OPERATIONAL COMPLEXITY.

The overall cost of deployment must be considered – RSVP-TE introduces a requirement for additional systems surrounding the network for both path placement and path monitoring. Additionally, costs are incurred based on operational training of staff where characteristics (e.g., reversion) differ to other protocol operation.



FOR PROTOCOL DESIGNERS:

Minimising the number of additional systems required <u>only</u> for the protocol minimises cost and deployment barriers – and makes for simplified roll-out.

LOOKING AT 5218: WAS RSVP-TE A SUCCESS?

Real problem?	AS DEMONSTRATED – EXPLICIT PATHS SEEM TO BE A REAL REQUIREMENT IN MULTI-SERVICE IP/MPLS NETWORKS.
No new hardware?	ONLY NEW CODE REQUIRED (IN GENERAL) – NEW HARDWARE IS A SIGNIFICANT CHALLENGE (LEGACY REMAINS!)
Existing ops/processes?	NO – SIGNIFICANT DIFFERENCE FOR REVERSION AND MONITORING – REQUIRES SPECIFIC TRAINING OF OPS STAFF.
Relieving operational pain?	NOT PARTICULARLY – BUT WE BALANCE OPERATIONAL COMPLEXITY AGAINST REDUCED PARALLEL DEPLOYMENTS.
Incrementally deployable?	CAN BE DONE FOR SUBSETS OF TRAFFIC – GOOD IN THIS RESPECT.
Good technical design?	INHERENT RELIANCE ON SOFT-STATE HAS SIGNIFICANT CHALLENGES – WAS THIS THE RIGHT CHOICE?

PARTIAL SUCCESS: WHEN THE BARRIER'S TOO HIGH.

RSVP-TE is deployed – but in more limited scenarios, with lower scale than envisaged – a partial success, limited by scalability and operational complexity.



Functionality

SEGMENT ROUTING/SPRING: CAN WE DO ANY BETTER?

Real world problems still exist – see Microsoft's SWAN or Google's B4 – systems implementing traffic engineering over and above those identified in this discussion.

CAN WE IMPLEMENT EXPLICIT PATHS WHICH CAN BE PERFORMANCE AWARE IN A MANNER WHICH SCALES TO TODAY'S REQUIREMENTS, AND LOWERS THE POINT AT WHICH THE COMPLEXITY BECOMES ACCEPTABLE?



WHAT IS SEGMENT ROUTING?



LABEL ADVERTISEMENT IN THE IGP.

FORWARDING BASED ON STACKED LABELS.

SEGMENT IDENTIFIERS.

Node A lo0: 172.16.12.1/32

Node-SID: 64

NODE SID: GLOBAL (INDEXED) LABEL ALLOCATION INDICATING SPT TO ADVERTISING NODE (SPECIAL PREFIX SID).

SEGMENT IDENTIFIERS.

Node B

Adj-SID: 1100

Node A lo0: 172.16.12.1/32

Node-SID: 64

ADJACENCY SID: LOCAL LABEL ALLOCATION INDICATING A LINK (OR SET OF LINKS) WITHIN THE IGP TOPOLOGY.



PREFIX SID: LOCAL LABEL ALLOCATION INDICATING AN IGP "LEAF" IP PREFIX (E.G, ATTACHED NODE).

Node-SID: 64



IGP-ANYCAST SID: GLOBAL LABEL ALLOCATION INDICATING REACHABILITY TO A CERTAIN RESOURCE OR FORWARDING PATH.

SPRING FORWARDING.

NODE-TO-NODE ALONG SPT:



NO NEED FOR LDP FOR FORWARDING TO NODES WITH NODE-SID, OR IP ROUTES WITH IGP-PREFIX-SID – CAN ELIMINATE LDP AND LDP-IGP SYNC.

SPRING TACTICAL TE: NODE-TO-NODE - TACTICAL TE:





SPRING TACTICAL TE: NODE-TO-NODE - TACTICAL TE:





SINGLE FORWARDING APPROACH FOR ALL FRR TYPES – NO ADDITIONAL CONTROL-PLANE REQUIRED.





IP FRR WITH SPRING.

EXPLICIT FORWARDING.



BASE COMPLEXITY: PATH CALCULATION.





RE-USE OF EXISTING CALCULATION MACHINERY: WHERE HEAD-END HAS VISIBILITY OF ALL REQUIRED ROUTING INFO. USE OF PATH COMPUTATION ELEMENT: GIVING HEAD-END ADDITIONAL VISIBILITY OR EXTERNAL INFO.

CHALLENGE: STACKED LABEL DEPTH.



DIFFICULT CHALLENGE – DEPENDENT ON OPERATOR TOPOLOGY AND VENDOR HARDWARE OPTIMISATIONS

COMPARISON TO RSVP-TE IN TERMS OF COMPLEXITY.



SOME CONCLUDING THOUGHTS...

OPERATOR'S ISSUES ARE WIDE AND VARIED...

IT'S EASY TO SEE PROBLEMS IN THE REAR-VIEW MIRROR...

CONSIDER THE COST OF CHANGE... Protocol design that looks to solve multiple sets of problems, for multiple operators are those that maximise their probability of success.

The collected thoughts in this presentation are based on issues in live operational networks – thinking about "what-ifs" in real networks is a good plan (but won't capture every issue).

Protocols that are already deployed have significant advantages over wholly new approaches – and may win out based on this – however, this <u>does not</u> prevent new deployments where there is sufficient value.



THANKS – QUESTIONS?

ROB.SHAKIR@BT.COM +44 (0)207 356 7378

ROB SHAKIR @ROBSHAKIR