

Multi-domain MPLS-TE

Latest developments in techniques for
computing inter-area and inter-domain
paths for traffic engineered MPLS

Adrian Farrel
CTO
Aria Networks Limited
adrian.farrel@aria-networks.com

Future-Net 2007



1



Agenda

- Ⓜ MPLS-TE Background
- Ⓜ What are Domains and Why Cross Them?
- Ⓜ Techniques for End-to-end Connectivity
- Ⓜ The Path Computation Element (PCE)
- Ⓜ Per-Domain Path Computation
- Ⓜ Crankback Routing
- Ⓜ TE Aggregation is bad!
- Ⓜ Backwards Recursive Path Computation
- Ⓜ Advanced Issues

2

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



MPLS-TE Background

- ④ MPLS-TE used to build “pipes”
 - ④ Direct traffic away from shortest paths
 - ④ Make best use of network resources
 - ④ Group traffic for common treatment
 - ④ Pseudowires, L3VPNs, scalability
 - ④ Quality guarantees through resource reservation
 - ④ Network repair and protection
 - ④ Fast Reroute (FRR)
 - ④ End-to-end protection
- ④ Signalled using RSVP-TE
- ④ Traffic Engineering Database (TED)
 - ④ Built from information distributed by the routing protocols
 - ④ Used to compute end-to-end paths

3

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Network Domains

“A domain is considered to be any collection of network elements within a common sphere of address management or path computational responsibility.” - RFC 4726

- ④ IGP areas
- ④ Autonomous Systems
- ④ Network layers
- ④ Client/server networks
- ④ Why cross domains?
 - ④ Because source and destination are not in the same domain!
 - ④ Multi-area and multi-AS networks, virtual POP, etc.
 - ④ Because one domain provides connectivity for another domain
 - ④ Client/server, multi-layer, VPN, etc.
- ④ How do we do it now?
 - ④ Manual stitching at domain boundaries
 - ④ Tunnel termination and reclassification of traffic at domain boundaries
 - ④ Careful off-line planning and management (e.g., FRR at domain borders)

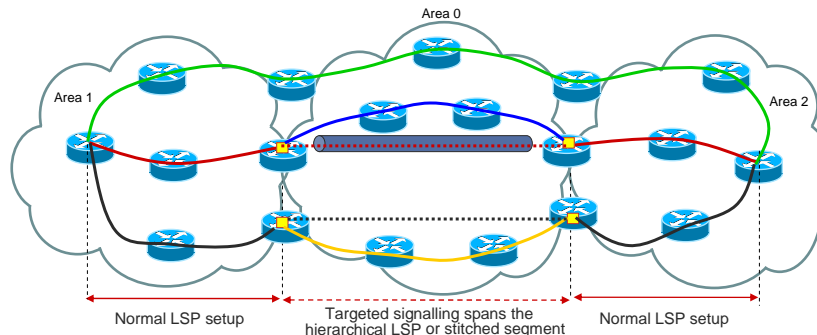
4

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Techniques for End-to-End Connectivity

- 1 Three techniques: contiguous, hierarchical, or stitched
- 2 Trade-offs
 - 3 Conceptual simplicity
 - 4 Administrative boundaries
 - 5 Data plane simplicity
 - 6 Reoptimisation and protection
- 3 Unanswered issues
 - 7 How to compute end-to-end paths
 - 8 How to select domain border nodes



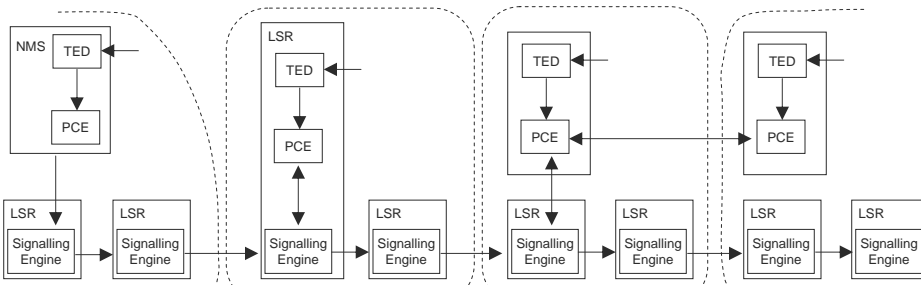
© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Path Computation Element

“An entity (component, application, or network node) that is capable of computing a network path or route based on a network graph and applying computational constraints” - RFC 4655

- 1 What's new?
 - 2 Nothing!
 - 3 A formalisation of the functional architecture
 - 4 The ability to perform path computation as a (remote) service

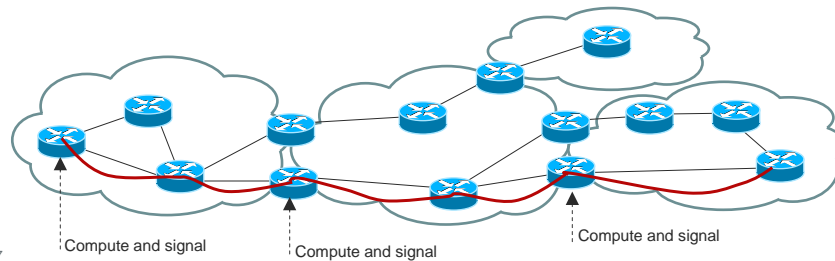


© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Per-Domain Path Computation

- Computational responsibility rests with domain entry point
- Path is computed across domain (or to destination)
- Works for contiguous, hierarchical, or stitched LSPs
- Which domain exit to choose for connectivity?
 - Follow IP routing? First approximation in IP/MPLS networks
 - Sequence of domains may be “known”
- Which domain exit to choose for optimality?

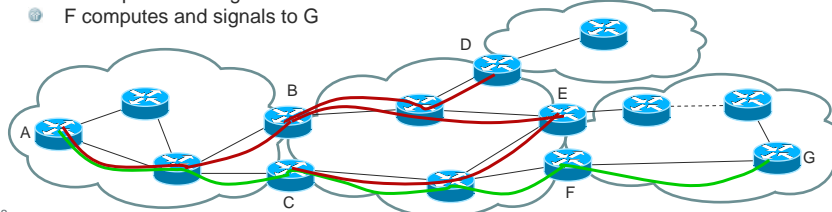


© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Crankback Routing

- A cure for connectivity, but not for optimality
- “Connectivity” means TE connectivity
 - May have IP connectivity, but insufficient resources
- May be painfully slow! “Informed random walk with wasted signalling”
 - A computes and signals to B
 - B computes and signals to D
 - D fails to compute and reports failure to B
 - B computes and signals to E
 - E computes to G, but no resources. Reports failure to B
 - B reports failure to A
 - A computes and signals to C
 - C computes and signals to E (can be avoided if E's previous report is passed around)
 - E computes to G, but no resources. Reports failure to C
 - C computes and signals to F
 - F computes and signals to G



© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.

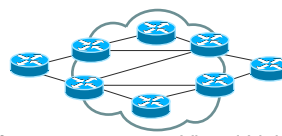


TE Aggregation is Not a Solution!

- The solution is “full TE visibility” but this does not scale
- TE aggregation looks very promising
 - Provide enough information to compute, but still scale
 - But aggregation reduces available information so optimality is in doubt
 - May hide connectivity issues
 - May cause confusing aggregation of information
 - May need frequent updates as internal information changes
- TE reachability also sounds good
 - Just provide information about which destinations can be reached
 - What does “reachability” actually mean?



Virtual Node aggregation
hides internal connectivity
issues



Virtual Link aggregation
needs compromises and
frequent updates

9

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Backward Recursive Path Computation

- PCE cooperation
 - Can achieve optimality without full visibility
 - “Crankback at computation time”
- Backward Recursive Path Computation is one mechanism
 - Assumes each PCE can compute any path across a domain
 - Assumes each PCE knows a PCE for the neighbouring domains
 - Assumes destination domain is known
- Start at the destination domain
 - Compute optimal path from each entry point
 - Pass the set of paths to the neighbouring PCEs
- At each PCE in turn
 - Compute the optimal paths from each entry point to each exit point
 - Build a tree of potential paths rooted at the destination
 - Prune out branches where there is no/inadequate reachability
- If the sequence of domains is “known” the procedure is neater

10

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.

BRPC Example

The diagram shows a network with three PCEs (PCE 1, PCE 2, PCE 3) and their domains of nodes:

- PCE 1 Domain:** Nodes A, B, C, D, E, F.
- PCE 2 Domain:** Nodes G, H, I, J, K, L, M, N, P.
- PCE 3 Domain:** Nodes Q, R, S, T, U, V.

Key connections include: A-B-C-D-E-F; G-H-I-J-K-L-M-N-P; Q-R-S-T-U-V; and inter-domain links: E-G, M-Q, P-Q, R-S, S-T, T-V.

- PCE 3 considers:
 - QTV cost 2; QTSRV cost 4
 - RSTV cost 3; RV cost 1
 - UV cost 1
- PCE 3 supplies PCE 2 with the tree
- PCE 2 considers
 - GMQ..V cost 4; GIJLNPR..V cost 7; GIJLNQ..V cost 8
 - HIJLNPR..V cost 7; HIGMQ..V cost 6; HIJLNQ..V cost 8
 - KNPR..V cost 4; KNPQ..V cost 5; KNLIJGMQ..V cost 9
- PCE 2 supplies PCE 1 with the tree
- PCE 1 considers
 - ABCDEG..V cost 9
 - AFH..V cost 8
- PCE 1 selects AFHIGMQTV cost 8

11

© 2005-2007. The Copyright in this presentation belongs to **Aria Networks Ltd.**

Advanced Computation Issues

- Inter-domain TE link information
 - For example, inter-AS links
 - Needs to be part of the information within a domain
- Path optimisation
 - Avoidance of “traps”
 - Trade-off of conflicting constraints
- FRR consideration during initial LSP placement
- Path diversity
 - End-to-end protection, load sharing, etc.
 - Link, node, domain, SRLG diversity
 - Avoidance of “traps”
- Reoptimisation
 - End-to-end or per-domain
 - “Shuffling” of deployed LSPs to free up stranded resources
 - May require migration strategies
- Different service types
 - Point-to-multipoint

12

© 2005-2007. The Copyright in this presentation belongs to **Aria Networks Ltd.**



The Future of Path Computation

- 1 Holistic Path Computation
 - 2 Solving the whole network is hard
 - 3 Balance conflicting constraints for different services
 - 3 Consider all services at once to avoid trap conditions
 - 3 Huge networks with thousands of services
 - 3 Needs to be adaptive to changes in topology and services
 - 3 Must be flexible to mixes of service types (P2P, P2MP, etc.)
 - 2 Necessary for full optimisation, but can it be achieved in real time?
- 2 Non-heuristic processes
 - 3 Conventional algorithms are deterministic and tuned to specific topologies and service types
 - 3 Non-heuristic processes can assess the whole network and all demands at once
 - 4 Can handle all topologies
 - 4 Can manage different service types
 - 4 Can trade-off conflicting constraints
 - 4 May produce a different, but correct solution each time
- 3 Highly sophisticated planning and modelling tools
 - 4 Multi-function
 - 5 Network failure analysis
 - 5 Capacity planning
 - 5 Rapid turn-around of network experiments
 - 5 Network re-optimisation
 - 4 Integrated planning and activation (NMS and PCE)
 - 4 On-line optimisation and reoptimisation
 - 5 Smart PCE
 - 5 Dynamic reconfiguration of networks with configured parameters, thresholds, and cost/risk/benefit analysis
- 3 Aria Networks Ltd. <http://www.aria-networks.com>

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.



Standardisation Status and References

- 1 RFC 4216: MPLS Inter-Autonomous System (AS) Traffic Engineering (TE) Requirements
- 1 RFC 4105: Requirements for Inter-Area MPLS Traffic Engineering
- 1 RFC 4726: A Framework for Inter-Domain Multiprotocol Label Switching Traffic Engineering
- 1 RFC 4655: A Path Computation Element (PCE)-Based Architecture
- 1 RFC 4206: Label Switched Paths (LSP) Hierarchy with Generalized Multi-Protocol Label Switching (GMPLS) Traffic Engineering (TE)
- 1 draft-ietf-ccamp-lsp-stitching: LSP Stitching with Generalized MPLS TE (work in progress)
- 1 draft-ietf-ccamp-inter-domain-pd-path-comp: A Per-domain path computation method for establishing Inter-domain Traffic Engineering (TE) Label Switched Paths (LSPs) (work in progress)
- 1 draft-ietf-pce-brpc: A Backward Recursive PCE-based Computation (BRPC) procedure to compute shortest inter-domain Traffic Engineering Label Switched Paths (work in progress)

14

© 2005-2007. The Copyright in this presentation belongs to Aria Networks Ltd.

Questions?

adrian.farrel@aria-networks.com

