

Network Operator Independent Resilient Overlay for Mission Critical Applications

Xian Zhang Chris Phillips

Networks Research Group, Department of Electronic Engineering and Computer Science
Queen Mary University of London
London, United Kingdom

xian.zhang@elec.qmul.ac.uk, chris.phillips@elec.qmul.ac.uk

Abstract—This paper proposes a Resilient Overlay for Mission Critical Applications (ROMCA); a novel operator-independent overlay architecture providing a resilient and reliable service across wide-area networks. Resilience is achieved by combining centralized topology construction control and distributed dynamic mapping of paths onto the overlay topology according to network conditions. ROMCA can mitigate the shortcomings of the underlying Internet network infrastructure and provide low recovery times in the event of network failure(s).

Keywords-ROMCA; resilience; reliability; mission critical

I. INTRODUCTION

Mission critical applications, e.g. remote control messaging, require high reliability and are very sensitive to network failure(s) and performance degradation, so there is considerable interest in developing a cost-effective scheme to provide the customers with a resilient delivery service for this kind of application. Currently, the Internet only provides “best effort” packet transport. Furthermore, the de facto inter-domain routing protocol - Border Gateway Protocol (BGP), used for routing across Autonomous Systems (AS) is characterized by re-convergence times of several minutes or longer [1]. These factors limit the ability of the Internet to support applications with higher reliability requirements; thus necessitating the exploitation of new strategies.

A novel architecture called ROMCA is proposed, which permits mission critical services to be set up and maintained despite uncertainty in underlying wide-area networks. Firstly, in Section II we provide a state-of-the-art review. Then, in Section III and IV, our architecture is illustrated and explained in detail and some typical operational examples are discussed. Finally, conclusions and ongoing work are briefly presented in Section V.

II. BACKGROUND

Overlay networks provide an effective way to overcome the shortcomings of the Internet in supporting new applications without introducing changes to the existing network layer [2, 3, 4, 5, 6, 7, 8 and 9]. For instance, Tapestry [3] and CAM-Chord [4] introduce an overlay layer to support peer-to-peer (P2P) and multicast applications on top of the Internet, respectively.

Till now, several significant architectures have been proposed aimed at improving Internet service performance, in terms of resilience, QoS (Quality of Service) or latency. RON [2] is the

first proposed overlay network aimed at improving the resilience of the Internet. The nodes in RON form a full mesh topology and use both active probing and passive measurements to monitor the Internet performance. If a working path undergoes failure(s) or it can no longer satisfy the prescribed performance requirement(s), the traffic is diverted to an intermediate RON node bypassing the unsatisfactory path. In contrast to the high convergence time of BGP, RON can achieve recovery times in the order of tens of seconds based on test-bed experiments. However, RON is application-specific and the number of overlay nodes is assumed to be no more than 50. It could incur scalability problems when widely deployed as RON utilizes a fully meshed topology and active probing and monitoring. The scalability issues of RON have been considered in [8], where hierarchically organized link state routing is employed.

The second type of overlay architecture is the one that is provider-dependent, whilst attempting to fulfil the QoS requirements of its customers. For example, the objective of SON [5] is to find overlay paths under certain bandwidth constraints. Another example is QRON [6]. It utilizes a kind of overlay node called Overlay Broker (OB) to construct a hierarchical topology, aimed at finding QoS-satisfied paths for end users. The nodes in QRON subscribe to an ISP for high bandwidth connections, thus the primary concern is how to provide the service in a cost-effective way. A moderate amount of work has been dedicated to discussing the overlay topology building process; one example is SIMON [7]. This employs a hierarchical distributed server mechanism to organize the intra-domain and inter-domain overlay nodes.

A third type of the architecture is exemplified by OHSR [12]. This is a simple, best-effort method for improving resilience proposed after extensive measurement research on the characteristics of the Internet path failure(s). OHSR tries to recover from path failures by routing indirectly using randomly chosen intermediaries without path monitoring. Exploiting the idea of OHSR, HORNS [13] proposes a heuristic node selection algorithm to support interactive real-time applications. The difference between the two is that HORNS improves OHSR performance by keeping the candidate node pool small and using end-to-end (ETE) delay as the selection criterion for the intermediary candidate.

ROMCA shares the same objectives as these overlay architectures, namely enhancing the performance of the Internet,

especially in supporting of applications with stringent requirements. However, its characteristics differ. It provides:

- **Network provider independence:** Based on the multi-domain Internet, ROMCA is designed to be a network-provider-independent overlay architecture, which can ensure the flexibility of its deployment across the wide area networks and avoid issues of inter-provider trust.
- **Scalable overlay topology:** In contrast to RON, the overlay nodes in ROMCA will be selectively picked from the Internet nodes and form into a partially meshed topology with layer-3 diversity. Moreover, it can be altered dynamically according to the results collected by exploiting ICMP-based traceroute methods. Thus, it can provide effective virtual linking between two overlay nodes that may be situated in different domains and performs path selection based upon customer service requests.
- **High Resilience:** Although, ROMCA shares the same objective of promoting the resilience of the Internet using alternative node(s) as RON and OHSR do, it employs a scalable topology that can guarantee the resilience of working and backup paths to some extent. Moreover, well-researched protection/restoration methods used in connection-oriented networks can be deployed in ROMCA overlay, which can provide the customer with highly resilient paths based on the dynamic overlay topology.

III. ROMCA ARCHITECTURE

As shown in Figure 1, ROMCA consists of a single *Overlay Directory Service (ODS)* and multiple *Overlay Gateways (OG)* chosen from different ASes.

The ODS is a centralized component responsible for service access and managing the overlay topology, including not only the acceptance and removal of overlay nodes, but also selection of the OG adjacencies. The ODS thus determines which virtual links will be set up between the OGs. However, it plays no part in the actual forwarding of traffic across the overlay. As [9] concludes, the overlay topology has an impact on overlay performance and physical network information is helpful in constructing the overlay topology. We therefore employ the ODS for organizing the overlay topology taking into account the position of potential overlay nodes. In terms of the IP path and topology discovery mechanism, [10] compares different traceroute methods and examines their performance. For example, ICMP-based traceroute tends to reach more destinations and collect “presence” information of a greater number of AS links as compared to other methods. Indeed, [11] confirms the superiority of traceroute method for discovering the Internet topology over alternative discovery methods. Thus, ICMP-based traceroute methods are used in supporting the construction of a layer-3 diversified overlay topology and are also used for performance monitoring of the virtual links between the OGs.

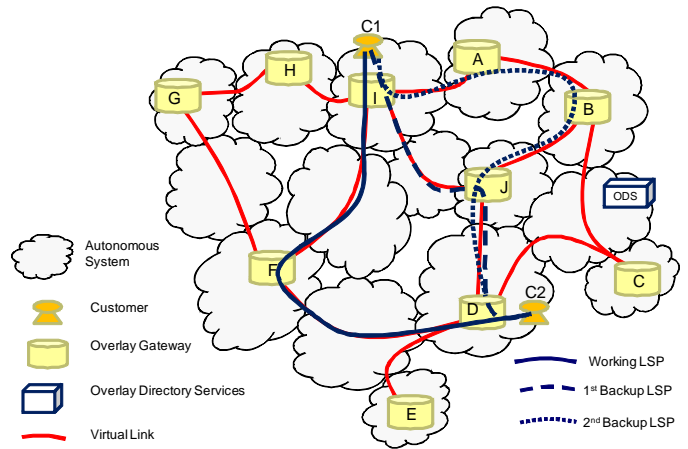


Figure 1: ROMCA Architecture

OGs are equipped with the following functions:

- Neighbour connectivity and performance information exchange between adjacent OG(s);
- Routing and performance information collection and dissemination in the overlay;
- Service provisioning, including establishing, maintaining and removing working and backup paths for customer traffic;
- Resilience-related functions, such as failure notifications to other OG nodes and the ODS;

Together the single ODS and multiple OG entities form the ROMCA architecture and are the means by which ROMCA provides resilience service to the end-users. The architecture itself is effectively hidden from the end-users, which simply know the public address of the ODS from which the appropriate OG points-of-presence can be ascertained. The example below explains the typical topology construction process of ROMCA.

It can be seen from the Figure 1 that the overlay topology is partially meshed and generally organized into inter-connected cycles, though stub connections are permitted. The virtual links between adjacent OGs nodes are chosen according to probing results and network performance measures. For instance, assume node G applies to the ODS to join the overlay. After retrieving the potential neighbouring information from the ODS and tracerouting to these corresponding nodes, G reports its findings to the ODS and is accepted into the ROMCA topology as it is a “valuable” transit node having Layer-3 diversified paths to H and F, i.e. from it multiple exit points from its local AS can be reached. There are also situations where stub nodes may get accepted according to their resilience and service access utility. For example, node C is viewed as a potential alternative path for B and D, so there are two virtual links connected from C. Whereas, the stub node E is only accepted into the overlay by establishing one virtual link connected to D. Node E is simply used as a service access point for the customer in its AS.

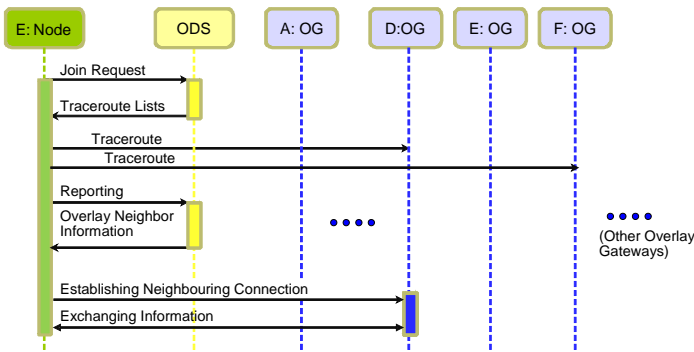


Figure 2: OG Node Joining Process

An example of the basic joining process for new nodes is depicted in Figure 2 (i.e. for node E) using a UML sequence diagram. As explained previously, the overlay topology is strategically constructed and maintained by the ODS, but a distributed mechanism for topology updating and network performance information flooding is needed so that OGs can maintain up-to-date performance information to enable them to efficiently establish working and backup paths for customer traffic. In our architecture, a flooding mechanism similar to that used in Open Shortest Path First (OSPF) and constraint-based routing protocol are deployed among the OGs. So, when the performance of a virtual link changes across a threshold, update packets will be flooded to all OGs so that they can store the updated information in their Link-State Database (LSD). This in turn will influence the routes chosen by the ingress OG for subsequent working and backup paths.

IV. OPERATIONAL EXAMPLES

A. Service Provisioning Example

To explain the service provisioning operation, consider the topology depicted in Figure 1, where Multi-Protocol Label Switching (MPLS) is used to define the working and protecting Label Switched Paths (LSPs). Consider a ROMCA customer (i.e. customer 1) located in the same Domain as OG I. The customer approaches the ODS providing the IP address of itself and the IP address of the destination customer (i.e. customer 2) from which the ODS can infer their proximity to the various OG nodes that are operational. In this case, the ODS knows the customer has no desire to become an OG; it simply wishes to exploit the overlay mesh to provide resilient pathways to a fellow customer in another AS. The ODS provides it with the nearest point-of-presence, i.e. the address of OG I, giving it a connection “ticket”. The ODS also informs the local OG, i.e. I, that it can expect an approach from customer A and the customer’s requirements as well.

When customer 1 contacts OG I, I checks the ticket details. In this case, the customer wishes to establish disjoint resilient paths to fellow customer 2. Using its LSD, OG I sends RSVP-TE path messages to OG D using as diverse links and nodes where possible. For example the working path may be taken to be via

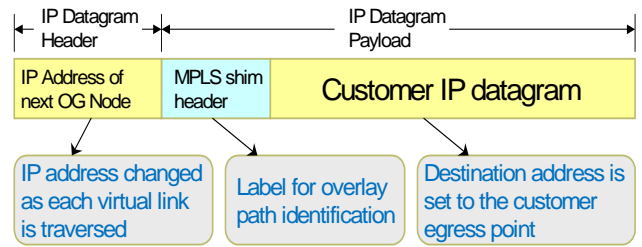


Figure 3: Packet Format used between OGs

OG I, F and D. The corresponding protection path could be: OG I, J and D. Once established, a FEC-to-Label binding entry is created at OG I, and customer 1 is informed that the service is ready. Traffic from customer 1 to 2 now uses IP to reach OG I (tunnelling IP in IP). There the packet will be encapsulated using the format depicted in figure 3.

Accordingly, the packet has an MPLS label pushed onto it and this is encapsulated in a datagram for OG F. At intermediate OGs, the MPLS shim layer is examined and the label is swapped and the traffic re-encapsulated and sent to the next-hop OG, and so on. At node D, the label is popped and the datagram delivered to customer 2 as per standard IP.

B. Resilience Implementation Scenarios

If a failure happens in a transit AS, it may take the ASes several minutes to re-converge and thereafter find the proper route to divert the traffic accordingly. But in ROMCA, as the neighbouring OGs exchange “hello” messages periodically (e.g. several seconds), the failure will result in loss of these heartbeat messages. Once the time threshold for neighbouring connectivity loss is reached, the OG(s) adjacent to the point(s) of failure will propagate the information over the virtual links to the ingress point(s), which can immediately update the FEC-to-Label binding so that the traffic is mapped onto the pre-configured diversified protection path(s). These paths avoid the failed AS and so can ensure that service delivery is quickly re-established.

As the service provisioning example shows, when ROMCA is used, traffic from customer 1 to 2 goes via the ingress OG I, and from there, follows a working path dictated by the LSP. Moreover, a protection path is also set up for resilience purposes. In the event of a failure in an AS lying between OG I and F, the ingress OG will switch the traffic to the protection LSP, i.e. to the path going from OG I, J to D, thus re-establish customer data packets transmission typically within seconds.

Another example is the dynamic mapping of LSPs according to the updated monitoring results. If the virtual link between I and J results in a longer delay than that of the path from I, via A, B to J, the backup LSP can be dynamically changed to the alternative backup LSP depicted in Figure 1, while the working LSP remains unchanged.

C. Performance Considerations

- Dynamic Overlay Topology

One feature of the ROMCA architecture is the adoption of a centralized means of topology construction, while providing services for customers in a distributed manner. The topology building strategy, utilizing the ICMP-based traceroute method among OGs, can ensure the creation of an efficient overlay topology and layer-3 diversity of virtual links to some extent, thus facilitating the resilience mechanisms deployed in the overlay. Moreover, as the performance and connectivity of the underlying paths change, the topology can be altered accordingly. It is also able to operate in a sparse deployment environment where multiple ASes may exist between adjacent Overlay Gateways, as is shown in figure 1.

As the failure of the ODS will affect the topology maintenance and service provisioning of new customers, backup strategies are possible such as mirroring the ODS in a similar manner to that employed by Open Shortest Path First (OSPF) for mirroring the Designated Router with a Backup Designated Router.

- High Resilience

Another feature of ROMCA is that both reactive and proactive methods can be exploited to meet the high resilience demands of the service customers. For instance, well-established protection mechanisms, such as 1:1, 1+1 and p-cycle, adopted in connection-oriented networks, can be employed. What's more, supplementary mechanisms such as machine learning can be incorporated to take advantage of the historical network performance data and make changes of working and backup paths before failure(s) occur. In the short term, failure(s) can be detected using regular "hello" messages exchanged between adjacent OGs; in the long term, we expect actions will be taken using prediction before outages interfere with customer services, by altering the virtual link arrangement.

- Low Recovery Time

Given each OG exchanges regular neighbour heartbeat messages (i.e. hellos) and the protection/restoration method adopted in overlay, we believe ROMCA can mitigate the slow convergence characteristic of BGP and achieve lower recovery times of the order of tens of seconds or even seconds on average. The absence of hellos along a virtual link triggers the switchover to backup paths (typically from the ingress OG), irrespective of the information being disseminated between ASes by BGP.

V. CONCLUSIONS AND ONGOING WORK

In this paper, an overlay architecture named ROMCA is introduced for supporting mission critical applications. Its mechanisms are explained with which the architecture can provide resilience with low recovery times in response to network

failure(s). In addition, ROMCA requires no specific support from network operators. This enables the overlay to be offered as a value-added enterprise service that can be deployed incrementally. It also capitalizes on the wealth of knowledge developed for resilience in existing circuit switched networks. Ongoing research is now evaluating its performance in terms of failure recovery ratio, overlay routing overhead and prediction performance both using a simulation platform and in field experiments.

REFERENCES

- [1] Craig Labovitz, Abha Ahuja, Abhijit Bose, Farnam Jhanian. "Delayed Internet routing convergence", *IEEE/ACM Trans. Networking*, Vol.9, No.3, pp. 293-306, June 2001.
- [2] David G. Anderson, Hari Balakrishnan, Frans kaashoek, Robert Morris. "Resilient Overlay Networks", *In Symposium on Operating Systems Principles*, pp. 131-145, 2001.
- [3] Ben Y. Zhao, Ling Huang, Jeremy Stribling, Sean C. Rhea, Anthony D. Joseph, John D. Kubiatowicz. "Tapestry: a resilient global-scale overlay for service deployment", *IEEE Journal on Selected Areas in Communications*, Vol. 22, Issue 1, pp. 41 – 53, Jan. 2004.
- [4] Zhan Zhang, Shigang Chen, Yibei Ling, Randy Chow. "Resilient Capacity-Aware Multicast Based on Overlay Networks", *IEEE International Conference on Distributed Computing Systems (ICDCS)*, pp. 565 – 574, 2005.
- [5] Zhenhai Duan, Zhi-Li Zhang, Yiwei Thomas Hou. "Service overlay networks: SLAs, QoS, and bandwidth provisioning", *IEEE/ACM Transactions on Networking*, Volume 11, Issue 6, pp. 870 – 883, Dec. 2003;
- [6] Zhi Li, Prasant Mohapatra. "QRON: QoS-aware routing in overlay networks", *IEEE Journal on Selected Areas in Communications*, Vol. 22, Issue 1, pp. 29 – 40, Jan. 2004.
- [7] Elaoud, M. McAuley, G. Kim, J. Chennikara-Varghese. "Self-initiated and self-maintained overlay networks (SIMONS) for enhancing military network capabilities", *Military Communications Conference, 2005. MILCOM 2005*. IEEE; Publication Date: 17-20 Oct. 2005 On Vol. 2. pp. 1147- 1151.
- [8] Hasegawa, G., Satoshi kamei, Masayuki Murata. "Emergency Communication Services Based on Overlay Networking Technologies", *Fourth International Conference on Networking and Services, ICNS 2008*. 2008.
- [9] Li Zhi, Prasant Mohapatra. "The impact of topology on overlay routing service" *Proceedings of IEEE INFOCOM 2004 - Conference on Computer Communications*, pp. 408-418, 2004.
- [10] Matthew Luckie et al. "Traceroute Probe Method and Forward IP Path Inference", *IMC' 08*, pp. 311-323, October 20-22, 2008.
- [11] Bin Yuan, Guoqiang Zhang, Yanjun Li, Guoqing Zhang, Zhongcheng Li. "Improving Chinese Internet's Resilience through Degree Rank Based Overlay Relays Placement", *ICC 2008*, pp. 5823-5827, 2008.
- [12] Gummedi, K.P., Harsha V. Madhyastha. "Improving the reliability of internet paths with one-hop source routing", *6th conference on Symposium on Operating Systems Design & Implementation*, Vol. 6. 2004.
- [13] Yin, C. et al. "Heuristic Relay Node Selection Algorithm for One-Hop Overlay Routing", *in Distributed Computing Systems Workshops (ICDCS)*, 2008.