

2007

Rerouting Technique for Faster Restoration of Preempted Calls

Iftekhar Ahmad
Edith Cowan University

Joarder Kamruzzaman

Daryoush Habibi
Edith Cowan University

10.1109/ICIS.2007.156 This article was originally published as: Ahmad, I., Kamruzzaman, J., & Habibi, D. (2007). Rerouting technique for faster restoration of preempted calls. Proceedings of ICIS 2007: 6th IEEE/ACIS International Conference. (pp. 676-681). Melbourne, Australia. IEEE Press.

Original article available [here](#)

This Conference Proceeding is posted at Research Online.

<http://ro.ecu.edu.au/ecuworks/1725>

Rerouting Technique for Faster Restoration of Preempted Calls

Iftekhar Ahmad#, Joarder Kamruzzaman* and Daryoush Habibi#

#School of Engineering and Mathematics, Edith Cowan University, Australia

*Gippsland School of Information Technology, Monash University, Australia

E-mail: {i.ahmad, d.habibi}@ecu.edu.au, Joarder.Kamruzzaman@infotech.monash.edu.au

Abstract

In a communication network where resources are shared between Instantaneous Request (IR) and Book-Ahead (BA) connections, activation of future BA connections causes preemption of many on-going IR connections upon resource scarcity. A solution to this problem is to reroute the preempted calls via alternative feasible paths, which often does not ensure acceptably low disruption of service. In this paper, a new rerouting strategy is proposed that uses the destination node to initiate the rerouting and thereby reduces the rerouting time, which ultimately improves the service disruption time. Simulations on a widely used network topology suggest that the proposed rerouting scheme achieves more successful rerouting rate with lower service disruption time, while not compromising other network performance metrics like utilization and call blocking rate.

1. Introduction

Resource reservation is the key technology that has gained increasing popularity as a research topic mainly due to the tremendous growth in *Multimedia and Distributed* applications that demand a predefined level of guarantee in terms of a number of parameters including end-to-end delay, packet loss rate, delay jitter and bandwidth availability. Depending on the application types and their dependency on resource availability, researchers [4], [5], [9] have proposed two types of resource reservation techniques: i) Instantaneous Request (IR) and ii) Book-Ahead (BA). IR calls are generally low bandwidth demanding calls and therefore reservations for these calls are attempted on the fly upon the arrival of requests. In contrast, BA calls generally demand high bandwidth and resource reservations for these calls are not often successful on the fly upon the arrival of the requests. For BA reservation, requests are required to be made in advance so that sufficient resources can be made available at the activation time of a BA application [4], [9]. Applications like multi-party video conferencing, video on demand, live TV broadcast programs,

telemedicine, grid computing and distributed simulations that require time critical start and demand high bandwidth are the potential candidates for BA reservations. Although both IR and BA reservation schemes are proposed to facilitate guaranteed QoS to the end applications, their co-existence at the same platform offers a number of key challenges and one such challenge is the need to preempt a number of on-going IR calls in order to supply the required resources for a BA call if resource scarcity arises at the starting time of that BA call.

Service continuity is a major element of a users' perceived QoS [1], [2], [3], [4], [15] and preemption of a connection in the middle of its lifetime causes a serious threat to its service continuity. In both wired and wireless networks, disruption of service continuity causes severe user dissatisfaction and long term revenue prospect of a network provider that depends heavily on user satisfaction is likely to suffer to a great extent. A number of strategies have been proposed in the literature that target to achieve low IR call preemption rate. Researchers suggested that IR call preemption rates can be successfully reduced at the routing and call admission control (CAC) stages. Ahmad *et al.* [4] presented a preemption-aware routing scheme that computes the preemption probability across various feasible paths and selects a path with the lowest chance of preemption for an in-coming IR call. Schelen *et al.* [5] suggested a look-ahead time based CAC scheme that reduced IR call preemption at the CAC stage. The motivation was to set aside resources for BA calls for a certain period in advance so that a BA call does not experience scarcity at the point of its activation. Ahmad *et al.* [6] improved this model by proposing an analytical method that determines the look-ahead time dynamically taking the changing traffic conditions into consideration. The strict partitioning of link capacity [7], [8] is another approach that divides the network resources into two disjoint subsets dedicated to each class of call and thereby eliminating the problem of preemption of any on-going calls. Optimizing the partition usage in such an approach is a major challenge and results in

drastically low network utilization [9]. Although all the aforementioned schemes are capable to successfully reduce the IR call preemption rate, it remains a daunting task to maintain it at a near zero level while maintaining satisfactory network performance in terms of other metrics like call blocking rate and resource utilization.

In cases where preemption of an IR call becomes unavoidable, the final remaining option is to initiate a rerouting process that makes an attempt to maintain the service continuity by restoring the connection through an alternative path that meets the demand of the connection. The most critical consideration for such a rerouting technique is the connection rerouting time, defined as the time duration required for reconnecting a connection through an alternative path once the connection along the primary path fails. Minimal service disruption time, which is the time interval during which reserved bandwidth for that connection is unavailable at any of the links across the path from the source to destination, can only be achieved if restoration of the connection can be made at the shortest possible time. A zero service disruption time is the ideal demand, but very hard to achieve due to delays involved in signaling of restoration messages and the need for reserving resources along the new path. A practical rerouting scheme therefore attempts to restrict the service disruption time to an acceptably low range so that the degradation in perceived QoS of the applications remains insignificant. A number of such rerouting schemes are discussed in the following section. However, none of the existing rerouting techniques known to the authors yields a service disruption time that is sufficiently low to reroute a connection without interruptions to the agreed QoS. This paper presents a new rerouting technique with the motivation to reduce the service disruption time and thereby provide improved QoS assurance.

2. Existing Rerouting Techniques

Researchers have classified the rerouting schemes mainly in two groups: i) proactive ii) reactive. In proactive schemes, resources are reserved and dedicated *a priori* along a back-up path for each connection so that the back-up paths can be used immediately following the failure of the primary connections. The proactive schemes include the *multiple copy* [10], the *dispersity routing* [11] and the *spare allocation* [12]. The proactive rerouting schemes provide lower service disruption time as there is no need to compute and reserve resources once the primary connection faces a problem. Maintaining a back-up path for each connection however, incurs

heavy costs and is infeasible in practice as it restricts future calls from using the resources, which results in an unacceptably low resource utilization and high call blocking rate. In the reactive schemes, an attempt to reserve resources and reroute the connection is initiated only after it is realized that resources allocated for a connection along the primary path is preempted. The reactive schemes are free from the overheads of maintaining resources along a back-up path, but causes long connection rerouting times and/or no rerouting at all, if the network is highly loaded. Doverspike [13] investigated the rerouting approaches from an implementation perspective and grouped them along three axes: link rerouting vs end-to-end rerouting; centralized vs distributed schemes; pre-computation vs dynamic computation rerouting schemes. The link rerouting vs end-to-end rerouting focuses on the network point (e.g., source or unlink node) at which rerouting has to be initiated, while the centralized vs distributed schemes concentrates on the controlling point of rerouting. In a centralized scheme, the rerouting process is controlled by a central point and hence suffers from the classical problem of unscalability, single point failure of the control unit, high latency and communication bottlenecks. The distributed approaches are free from these problems, but provide sub-optimal solutions. The time it takes to select an alternative path is the point of focus for the pre-computation vs dynamic computation rerouting schemes. Banerjee *et al.* [14] investigated rerouting schemes along three key components: locus to reroute, reroute timing and retry model. The 'locus to reroute' model concentrates on the link (local) vs end-to-end (global) choice of rerouting schemes. The immediate upstream node of the failed link/node takes the responsibility to determine a path segment from that node to the destination in a local rerouting model and thereby reroutes the connection. In the 'end-to-end' rerouting model, any failure information is sent to the source node and the source node determines a path from that the source to the destination. The reroute component determines the start time for the rerouting attempt and the decision concerning this time is governed by the reroute timing model. Immediate, random and sequential are the possible approaches for reroute timing. Immediate timing is the least cooperative approach and needs an instant solution. Immediate timing initiates the rerouting process as soon as the failure information is reported. The reroute time is determined by generating a random value from a uniform distribution over an interval of time for random timing approach. This approach provides some levels of cooperation among the nodes trying to

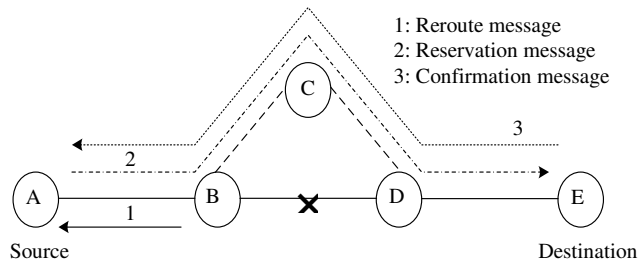


Fig. 1: Connection rerouting in reactive approach.

reroute. In sequential timing, all rerouting attempts are made sequentially with only one controlling node initiating a rerouting attempt at any given time. The 'retry model' controls the number and timing of retry attempts if the attempt for connection restoration initially fails. Immediate retry tries to reroute the connection immediately following the failure of the previous attempt, while a delayed retry attempts to reroute after a randomly selected waiting interval.

While proactive and reactive approaches have been the major groups for rerouting techniques, Ahmad *et al.* [18] recently proposed a variation of the conventional approaches, which is particularly suitable for a communication network where resources are shared between BA and IR calls. The motivation was to estimate the resource scarcity at the BA activation points and initiate the rerouting of IR calls in advance so that the calls that are likely to be preempted are rerouted before experiencing the actual preemption. The results were found very promising in terms of both the successful rerouting rate and service disruption time. In this paper, we propose an improvement of our previous work by focusing on how to improve the rerouting time required for reserving the resources across the newly selected path.

3. Rerouting Scheme for Faster Recovery

As emphasized in the previous section, proactive schemes are not suitable for a commercial network due to low resource utilization. Since high resource utilization is a major consideration in a commercial network, reactive rerouting is more attractive, which forms the basis of the proposed scheme. The connection rerouting time in reactive rerouting scheme is comprised of delays at three stages: i) time required for the rerouting message to reach the source/uplink node (Step 1 in Fig. 1), ii) Time required for source/uplink node to compute an alternate path, and iii) Time required for the reservation message to traverse the whole round trip path to reserve the requisite bandwidth (Steps 2 and 3).

The time required for the rerouting message to reach the source/uplink node depends on the choice of the locus of rerouting, which determines whether the uplink or source node will attempt to reroute the connection. Banerjea *et al.* [14] suggested that source (global) rerouting always outperforms uplink (local) rerouting in terms of successful rerouting rate and since unsuccessful attempts of rerouting result in complete termination of the service, the choice of locus that contributes to higher successful rerouting rates should be exercised in QoS-enabled networks. The source rerouting requires a time duration to select an alternate path from the source to the destination and the length of this duration depends on both the routing algorithm and its computational complexity. Once a path has been selected by the source rerouting, standard resource reservation protocols send a reservation message across the new path which cross-checks (CAC) the feasibility of allocating resources to that connection and reserves the necessary resources given that the CAC process succeeds across the new path. This involves a complete round trip of the reservation message. If the round trip message reaches the source with the confirmation of resource being reserved, the source starts to transmit data via the new path. In such cases, the service disruption time equals the connection rerouting time. In this paper, we propose a technique to reduce the service disruption time by initiating the rerouting process from the destination node instead of the source node.

Let us consider a network scenario as shown in Fig. 2. Let t_s be the nearest BA call activation time (potential resource scarcity point), T_f the time required for the rerouting message to reach the locus of rerouting, T_p the time required by the routing algorithm to select a path from the locus to the destination, and T_r the time for a reservation message to make a round trip travel and reserve bandwidth across the path. The total rerouting time T_R can then be expressed as:

$$T_R = T_f + T_p + T_r \quad (1)$$

If n_1 is the total number of links from the failure node to source node and average delay for message traversal between two nodes is Δ_1 , then T_f is expressed as:

$$T_f = n_1 \Delta_1 \quad (2)$$

If n_2 is the length of newly selected path and the average delay for call admission control per link is Δ_2 , then T_r becomes:

$$T_r = n_2 (\Delta_1 + \Delta_2) + n_2 \Delta_1 \quad (3)$$

Using Eq. (2) and (3), Eq. (1) takes the following form

$$T_R = n_1 \Delta_1 + \Omega + n_2 (\Delta_1 + \Delta_2) + n_2 \Delta_1. \quad (4)$$

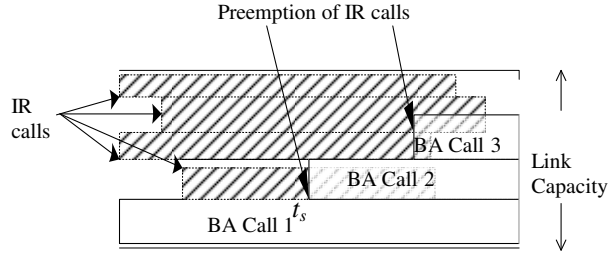


Fig. 2: Preemption scenario of IR calls.

where $T_p = \Omega$ represents the time required by the routing algorithm to select a path from the source to destination. For a network with E edges and V nodes, the maximum length that a link failure message may have to traverse is $(|E|-1)$ in source rerouting. The length of the alternative path in the worst case may contain all the links except the failed one. Therefore, the reservation message in the worst case may have to traverse up to $(|E|-1)$ links. The routing algorithm (Dijkstra's shortest path algorithm [17]) in such case will require the time complexity of $O(n \log n)$ where n equals $(|E|-1)$. For a network with a known system performance (hardware and software), the worst case routing complexity represents a value (Ω) in real time. Using the above notations Δ_1 and Δ_2 for delay in message traversal and CAC, respectively, the worst case restoration time T_{WR} can be expressed as

$$T_{WR} = (|E|-1)\Delta_1 + \Omega + (|E|-1)(\Delta_1 + \Delta_2) + (|E|-1)\Delta_1 \quad (5)$$

While the above formula can be generically used for all types of network, the calculation of maximum rerouting time can be made more realistically once the connectivity information is available. For a network with fixed and known physical connectivity, it is possible to compute the longest possible path in the network between any pair of nodes using all pair longest path algorithm which is essentially the all pair shortest path algorithm with a modification in the objective function. Denoting N as the length of all pair longest path, the worst case rerouting time with known connectivity takes the following form

$$T_{WR} = (|N|-1)\Delta_1 + \Omega + (|N|-1)(\Delta_1 + \Delta_2) + (|N|-1)\Delta_1. \quad (6)$$

3.1 Proposed Destination Driven Rerouting

Both forward pass and reverse pass [14] reservation schemes require the reservation message to traverse a round trip path from the source to destination as illustrated in steps 2 and 3 in Fig. 1. This is particularly important during the call connection set-up process because the connection request is made at the source node and the source node must be aware of the status of connectivity before it starts to send data. The routing process is executed at the source node and a round trip

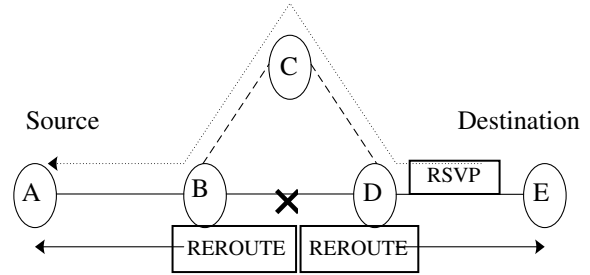


Fig. 3: Destination Driven Reservation for rerouting.

message that carries a confirmed status about connectivity is required by the source node before admitting the connection.

As indicated in Eq. (4 & 6), the round trip time for reservation message is one of the sources of delay in connection rerouting time, specially in a large size network. Reduction of the traversal delay for resource reservation message will decrease the connection rerouting time. While a round trip traversal of the reservation message is a necessity during the initial call setup process, it can be relaxed to a one-way trip during connection rerouting stage. Once a connection rerouting process is initiated, the uplink node (B) sends a failure message to the source node (B->A) as shown in Fig. 3. At the same time the downlink node (D) also sends a rerouting message to the destination node (D->E). When the source node receives the message, it may stop transmission or continue to overflow the system depending on the traffic engineering policy. The destination node upon receiving the message attempts to establish a path to the source node. This has to be done by executing the routing algorithm at the destination node followed by sending reservation message in the forward pass technique (E->D->C->B->A). If the reservation message reaches the source node with success at each node along the path, the source node immediately starts to send data along the new path which completes the whole rerouting process. For destination driven rerouting scheme, the reservation message is required to travel from destination to source only. This saves the source to destination traversal time of reservation message required in existing source driven reservation scheme. Mathematically, the saving in time is $n_2\Delta_1$ where n_2 is the length of the restored path. The proposed destination driven forward pass reservation scheme requires the following connection rerouting time

$$T_R = n_1 \Delta_1 + \Omega + n_2 (\Delta_1 + \Delta_2). \quad (7)$$

Following the same analysis as done earlier, for a network with known connectivity, the maximum rerouting T_{DWR} in destination driven reservation technique reduces to

$$T_{DWR} = (|N|-1)\Delta_1 + \Omega + (|N|-1)(\Delta_1 + \Delta_2). \quad (8)$$

The difference between Eq. (8) and (6) therefore indicates the improvement in numerical term achieved by the proposed scheme.

4. Simulation Results

The topology that has been used for the simulation represents a typical ISP network that follows the ATT backbone network structure and has been simulated in previous studies [4], [16], [17]. Bandwidth demand of each BA and IR call is uniformly distributed in the range of 1 to 2 Mbps and 64 to 256 kbps, respectively. Lifetime of BA and IR calls is exponentially distributed with a mean of 300s and 90s, respectively. Arrival of BA and IR calls is assumed to follow a Poisson distribution with a mean arrival interval of 10s and 200ms, respectively. Average propagation latency for each link is considered as 10ms and average time requirement for CAC is considered as 2ms for each link [14]. Since the study is based on BA reservation, each simulation is repeated for different BA limits β . BA limit β sets the normalized limit on link capacity that the aggregate BA load can use so that starvation for IR load is avoided. In our simulations, the average service disruption time per successfully rerouted call, successful rerouting rate (SRR) and SRR with tolerable service disruption were investigated. Successful rerouting rate is the ratio of the total number of successfully rerouted calls to the total number of attempted rerouting calls, with the former including those preempted calls that are reconnected through alternate paths at zero or finite time disruption in service continuity. We also investigated the effect of the proposed technique on network utilization and call blocking rates. We implemented the proposed destination driven rerouting scheme on top of the rerouting in advance scheme [18] (i.e., for simulation of the proposed scheme, rerouting process is initiated in advance, but reservation across the new path is initiated from the destination node instead of the source node). We compared the performance of the proposed destination driven rerouting in advance (DRA) scheme against a standard reactive rerouting (SR) scheme and a recently proposed source driven rerouting in advance (RA) scheme [18].

Figure 4 shows the average service disruption time per successfully rerouted IR connection in the proposed DRA and existing SR and RA schemes. The results indicate that the proposed rerouting scheme achieves the lowest average service disruption time, so validating the benefit secured of destination driven reservation. The proposed DRA scheme consistently

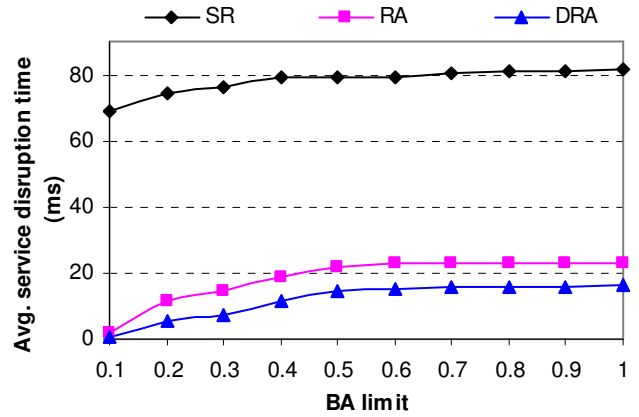


Fig. 4: Average service disruption time across different BA limits.

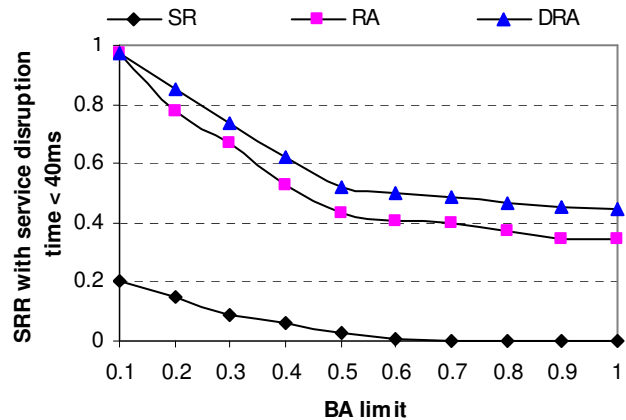


Fig. 5: Successful rerouting rate across different BA limits.

outperforms the SR and RA scheme by a margin of up to 66ms and 8ms, respectively. The average service disruption time increases with increasing BA limits because the net network load increases due to the increasing offering of BA loads. With the increased load, the length of feasible alternative paths is often longer and hence it takes longer time to reroute IR calls.

Figure 5 shows the successful IR call rerouting rate when the maximum allowable service disruption time is restricted to a limit of 40ms. Here, the limit assumes that if the connection can be rerouted within 40ms limit following the preemption, the degradation in perceived QoS will be still tolerable to the users as per the service level agreement. The proposed rerouting scheme outperforms the SR and RA scheme by a margin up to 45% and 10%, respectively. This is a promising improvement as many users in such cases will remain satisfied even if their connections were preempted at some stages. This satisfaction may prove significant for the revenue prospect of the network provider in the long run. The amount of data loss observed in different rerouting schemes is depicted in Fig. 6. The figure

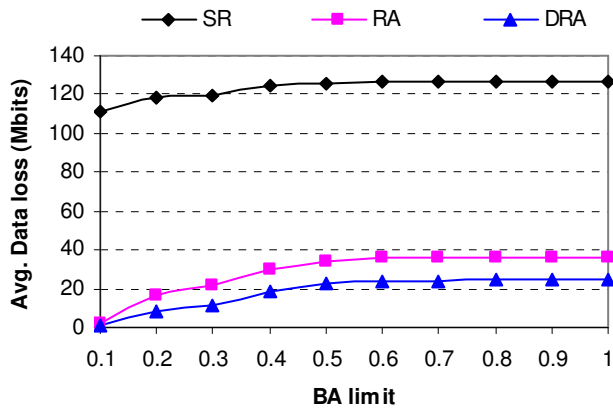


Fig. 6: Average data loss across different BA limits.

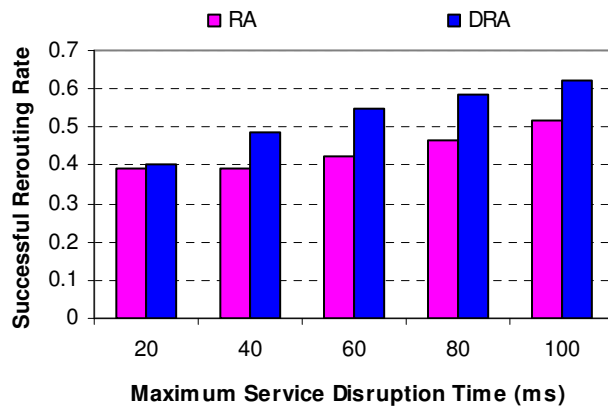


Fig. 7: Successful rerouting rate for various time limits.

confirms that the proposed DRA scheme achieves the lowest data loss across different rerouting schemes. Lower data loss rate has the advantages of yielding higher net utilization of the resources and higher revenue earning. Figure 7 shows the successful rerouting rate for various maximum allowable service disruption times at a BA limit of 0.7. The figure suggests that the improvement achieved by the proposed rerouting scheme is consistent across different values of maximum allowed service disruption time. The importance of DRA is increasingly realized for higher values of maximum service disruption time. We also observed the call blocking rate in different rerouting schemes and the proposed DRA scheme performed comparably with the RA scheme.

5. Conclusion

In this paper, an improvement of a recently introduced rerouting in advance scheme has been proposed with the objective to improve the service disruption time and successful rerouting rate of preempted IR calls when both IR and BA calls share the same resources. In this paper, the time required for reservation messaging delay is improved by engaging the destination node to

initiate the reservation process instead of the source node. The amount of time that the proposed scheme can improve is indicated mathematically in this paper. While the benefit of destination driven rerouting is evident for all types of network scenarios, the proposed scheme yields the most promising results when the failure node is close to the destination node and the path length is long.

6. References

- [1] W. C. Hardy, *QoS measurement and evaluation of telecommunications quality of service*: John Wiley & sons Ltd, New York, 2001.
- [2] M. Campanella, P. Chivalier, and N. Simar, "Quality of Service Definition," <http://www.dante.net/sequin/QoS-def-Apr01.pdf>.
- [3] B. Awerbuch, Y. Azar, and S. Plotkin, "Throughput-competitive on-line routing," *Proc. IEEE Annual Symposium on Foundations of Computer Science*, pp. 32-40, 1993.
- [4] I. Ahmad, J. Kamruzzaman, and S. Aswathanarayanan, "Preemption-Aware Routing for QoS-Enabled Networks," *Proc. IEEE Global Telecommunications Conference (GLOBECOM 2005)*, USA, 2005.
- [5] O. Schelen and S. Pink, "Sharing resources in advance reservation agents," *Journal of High Speed Networks*, vol. 7, no. 3-4, pp. 213-218, 1998..
- [6] I. Ahmad, J. Kamruzzaman, and S. Aswathanarayanan, "A dynamic approach to reduce preemption in book-ahead reservation in QoS-enabled networks," *Computer Communications*, vol. 29, no. 1, pp. 1443-1457, 2006.
- [7] K. W. Ross, *Multiservice Loss Models for Broadband Telecommunication Networks*: UK: Springer-Verlag, 1995.
- [8] D. Ferrari, A. Gupta, and G. Ventre, "Distributed advance reservation of real-time connections," *Proc. NOSSDAV*, pp. 15-26, 1995.
- [9] A. G. Greenberg, R. Srikant, and W. Whitt, "Resource sharing for book-ahead and instantaneous-request calls," *IEEE/ACM Trans. Networking*, vol. 7, pp. 10-22, 1999.
- [10] P. Ramanathan and K. G. Shin, "Delivery of time-critical messages using a multiple copy approach" *ACM Transaction Computer Systems* vol. 10, no. 2, pp. 144-166, 1992.
- [11] A. Banerjee, "Simulation study of the capacity effects of dispersity routing for fault tolerant real-time channels" *Proc. ACM SIGCOMM*, pp. 192-205, 1996.
- [12] S. Han and K. G. Shin, "Efficient spare resource allocation for fast restoration of real-time channels from network component failures" *Proc. IEEE RTSS*, pp. 99-108, 1997.
- [13] R. Doverspike, "A multi-layered model for survivability in intra-LATA transport networks," *Proc. IEEE GLOBECOM*, pp. 2025-2031, 1991.
- [14] A. Banerjee, "Fault recovery for guaranteed performance communications connections" *IEEE/ACM Trans. Networking*, vol. 7, no. 5, pp. 653-668, 1999.
- [15] I. Ahmad, J. Kamruzzaman, and S. Aswathanarayanan, "Preemption Policy in QoS-Enabled Networks: A Customer Centric Approach", *Journal of Research and Practice in Information Technology*, Australia, 2007.
- [16] G. Banerjee and D. Sidhu, "Comparative analysis of path computation techniques for MPLS traffic engineering," *Computer Networks*, vol. 40, no. 1, pp. 149-165, 2002.
- [17] T. H. Cormen, C. E. Leiserson and R. L. Rivest, *Introduction to Algorithms*, MIT Press, USA.
- [18] I. Ahmad, J. Kamruzzaman and L. Dooley, "Look-Ahead rerouting of preempted calls," *Proc. ATNAC 2006*.