

# A Policy-Based Bandwidth Resource Provisioning Architecture

David Chieng and Alan Marshall

*School of Electrical & Electronic Engineering, The Queen's University of Belfast, Ashby Bld, Stranmillis Road, BT9 5AH Belfast, Northern Ireland, UK.*

**Abstract:** We present a policy-based, agent-enhanced resource provisioning architecture, which facilitates flexible and quantitative end-to-end bandwidth reservation and management on a per-user, per-application or per-flow basis. Based on this architecture, various session level Service Level Agreement (SLA) negotiation schemes involving bandwidth allocation, service start time, guaranteed session duration can be introduced. The results show how these negotiation schemes can be utilised for the benefits of both network users and providers.

**Keywords:** Quality of Service, Agents, Bandwidth Reservation, Bandwidth/QoS Negotiation, Service Level Agreement, Policy System,

## 1. INTRODUCTION

Current IP network infrastructures are undergoing rapid transformations, i.e. from providing mere connectivity, to a wider range of flexible and complex tangible services involving Quality of Service (QoS). With QoS support, it is now possible to deliver different levels of service quality for a network application through Service Level Agreement (SLA). However, in today's already complex network environment, provisioning such services is a great challenge. Firstly, service and network providers will have to deal with a myriad of user requests that come with diverse QoS or SLA requirements. The providers will then need to make sure that these requirements can be delivered accordingly. Matching these service requirements to a set of control mechanisms in a consistent manner remains an area of weakness within the existing IP QoS architectures [1]. As well as

---

The original version of this chapter was revised: The copyright line was incorrect. This has been corrected. The Erratum to this chapter is available at DOI: [10.1007/978-0-387-35620-4\\_43](https://doi.org/10.1007/978-0-387-35620-4_43)

D. Gaïti et al. (eds.), *Network Control and Engineering for QoS, Security and Mobility*  
© IFIP International Federation for Information Processing 2003

the need to provide dynamic traffic management, monitoring and control, other operation or policy issues such as service management, security, customer management, accounting and billing also must also be addressed.

To address the issues above, we propose a novel agent-enhanced Bandwidth Resource Provisioning Architecture (BRPA), which facilitates dynamic, flexible and quantitative end-to-end bandwidth resource provisioning and management on a per-user, per-application or per-flow basis. The inherent characteristics of agents such as autonomy, adaptability, social abilities, etc offer many advantages in this dynamic, complex, distributed and heterogeneous network environment [2]. For example, a service provider agent can play a major role in guiding and deciphering users' requests, and is also able to respond quickly and effectively to their requests. The tasks also include service brokering, QoS specification and reconfiguration, pricing negotiation, etc. This is important as many end users/customers and service or network providers in general are still unable to specify SLAs in a way that benefits both parties.

This paper is organised as follows: Section 2 first introduces the BRPA's architecture. Section 3 presents the SLA negotiation schemes and analyses. Finally, the related work and the conclusions are presented in section 4 and 5 respectively.

## **2. BANDWIDTH RESOURCE PROVISIONING ARCHITECTURE (BRPA)**

To provision QoS, some form of resource reservation or allocation (i.e. bandwidth) need to be applied either implicitly or explicitly. This parameter is emphasized since it is the single most important factor that affects the QoS. The limits for delay, jitter and buffer size can be determined by the bandwidth reserved for a flow [3]. Most existing resource allocation mechanisms such as Resource Reservation Protocol (RSVP) [4] and ST-2 [5] only consider initial availability and do not take into account changes in future availability. BRPA not only provides immediate reservation, it is also able to reserve resources in advance. The architecture consists three sets of databases and a Resource Manager (RM). These are described as follows:

### **a) Resource Reservation Table (RRT)**

This database as illustrated in figure 1 comes in the form of a resource/time table so that network provider can lookup and allocate network resources (bandwidth) at present and also in the future. The RRT tells the network provider about the current or future available and reserved bandwidth resource of any link, at anytime. Here, we can see that at time  $t$ ,

some resources in the future have been reserved. A minimum unit or the granularity of 'reservable' session  $\tau$  must be defined. This can be 5 minutes, 1 minute or one second depending on the network provider's policy. Similarly, the minimum unit of 'reservable' bandwidth is defined as  $b$ . It can be quantified in units of kbps, 10kbps, etc. The RRT only indicates the reservation load and does not represent the actual link utilisation due to real-time traffic flows. It can be observed from figure 1 that request  $R_x$ 's guaranteed bandwidth request cannot be honoured throughout the entire requested duration  $T_x$ . In this case the requester has a number of options. He or she can either reduce the amount of bandwidth required, postpone the reservation session start time, cut short the session duration or accept the compromise that for a short period of time, their bandwidth will drop off.

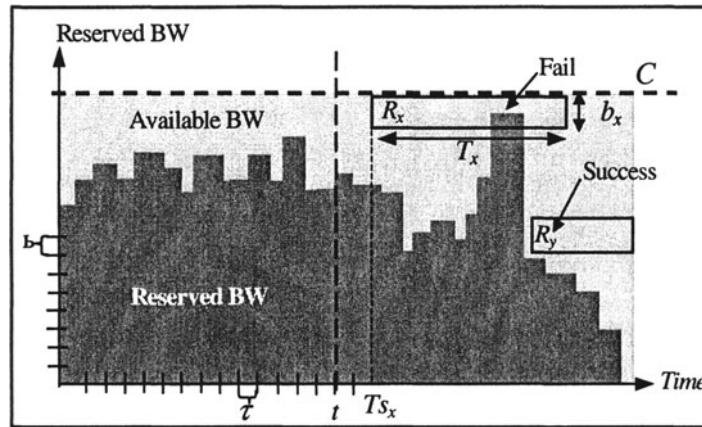


Figure 1. Resource Reservation Table and Reservation Processes

**b) Path Tables (PT)**

The path table as shown in figure 2 is introduced besides the native short path first routing table. It provides a distinct route or path info from a source node to a destination node and vice versa on a hop-by-hop basis.

		Destination			
		Node ID	0	1	2
Source	0	-	A2B2	A1B3D3	A4C3
	1	B1A1	-	...	...
	2	...	...	...	...
	3	...	...	...	...
	3	...	...	...	...

Figure 2. Path Table

For example from end node 0 to node 1, the traffic needs to transverse via router A output port 2 and then router B output port 2 (A2B2). If there is more than one end-to-end path for a source-destination pair, multiple path

tables need to be created. The PTs are similar to MPLS' Traffic Engineering Link-State Database (TE-LSDB) [6] that is used to provision distinct paths with QoS preferences. The combination of RRTs and PTs enables end-to-end resource reservation at present and also in the future. The Resource Manager (RM) is employed to manage all these resources. At the link level, an agent may be assigned to manage each RRT.

### c) User SLA Database (USD)

This database stores individual customer's SLA information. Each element stores a user's service ID, the respective bandwidth allocation ( $b$ ), agreed service activation time ( $T_s$ ), session duration ( $T$ ), end-to-end path ID or routing option, billing option, and also other rules and policies. These records are constantly monitored and updated. With this record, the network provider is able to manage the connection on a per user basis. The value added services are billed and delivered on a per-transaction or per-session basis, which can be based on bandwidth usage, or flat-rate pricing scheme. When the duration of resource reservation  $T$  expires, the connection will automatically revert to best effort mode. Alternatively, the allocated bandwidth will sustain until it is needed by other incoming non-preemptable sessions. In other words, when the non-preemptable session expires, it will automatically switch to preemptable mode. The user may need to re-negotiate in order to extend the reservation session. This parameter is typically used for Video on Demand (VoD) and news broadcast type services where service sessions are known in advance. In a situation where the desired session time cannot be specified or it is not known a priori such as in IP voice calls, an alternative scheme is necessary. For example these services can be automatically assigned to preemptable mode.

## 2.1 Overall Operation and Benefits

Upon receiving an SLA request from the UA, the RM 1) first retrieves the path information from PTs; 2) it then checks the resource availability at the RRTs. Normally the shortest available path will be quoted first. 3) When all the SLA terms and conditions are satisfied, the information will be stored in the USD. 4) The RM then propagates the resource requirements to the RAs, which configure the routers accordingly. To increase resiliency, a backup path can also be provisioned together with the primary path.

The admission control policy in BRPA is contrary to IntServ's RSVP [4] peer-to-peer approach. Using RSVP, routers normally decide locally whether to accept or reject resource reservation. This is done on a hop-by-hop basis along the shortest-path tree as RSVP is designed to work with existing routing protocols. Hence, local routers are not aware of any alternative paths

available in case of a reservation failure. In the BRPA architecture, the alternative path can be allocated via BRPA if the shortest path is saturated. However, RSVP can still be applied here if the DM first makes sure that the required resources are available along the shortest path before the RSVP RESV message is allowed to propagate. The proposed architecture also does not add to scalability problems since a path's info and link's states are stored in a high-capacity centralized domain server i.e. domain/resource manager server and not in the routers themselves.

The use of agents also enables dynamic resource negotiations and therefore promotes greater flexibility. In conventional reservation scheme, reservation will fail even if only 1% of the requested bandwidth cannot be met. In the proposed architecture, an alternative can be resolved through negotiations. Furthermore as mentioned in [7], the RSVP-based reservation is not suitable for all applications especially when sources define the resource requirements. Here, the agents can be employed to optimise resource usage bilaterally. In this architecture, autonomous agents also play an important role in enhancing service discovery process i.e. via advertisement. This is essential in heterogeneous network environments as not all networks may offer services with QoS options.

### 3. SLA NEGOTIATION SCHEMES EVALUATION

With flexible resource provisioning facility enabled via BRPA, dynamic SLA negotiations can take place. This event can be exploited for the benefits of both negotiating parties. For example, by specifying tolerances or alternative options in an SLA request, the chance of having it being rejected can be reduced. The User Agent (UA) can be employed to carry out strategic negotiation tasks such as getting the highest individual SLA optimisation in terms of QoS and price [8]. The Network Provider Agent (NPA) on the other side can facilitate service customisation, and at the same time trying to maximize the company's interests such as revenue. However, due to the limited space only two session-level negotiations schemes involving bandwidth i.e. Bandwidth Negotiation at Resource Limit (BNRL) and Bandwidth Negotiation at Predefined Load Level (BNPLL) are presented. The study focuses on the following three performance issues: 1) *Request Rejection Probability* ( $p_{rej}$ ), 2) *Percentage Mean Utilisation or Mean Reservation Load* ( $\%R_r$ ) and 3) *Bandwidth Satisfaction Index* (BWI). BWI defined as the mean ratio of bandwidth granted over bandwidth requested or  $BWI = \text{mean}(b_g/b_r)$  for all the users. It may be viewed as corresponding to the overall users' satisfaction.

### 3.1 Simulation Parameters

The followings are the simulation parameters used: The mean session duration requested by a user,  $T$  is fixed at 300s. The smallest session unit or reservation duration granularity,  $\tau=1$ s. The bandwidth unit requested per call,  $b_r$ , ranges from 1 unit to 156 units, is generated according to random exponential mean distribution. The reservation bandwidth granularity,  $b=1$  unit. This is equivalent to a range of 64kbps up to 10Mbps per call if 1 unit = 64kbps. The link capacity,  $C$  is fixed at 1562 bandwidth units or equivalent to 100Mbps if  $b=64$ kbps. This emulates the demand for different multimedia services such as voice, hi-fi music, video telephony, up to high-bandwidth video conferencing or VoD sessions with different quality levels and hence bandwidth requirements. In this study, only a unidirectional bandwidth reservation across a single link is considered although the same utility can be applied for both uplink and downlink. For an advance reservation (AR) request, the minimum session start time,  $T_{s_{min}}$  is fixed at 600s or 10 minutes from the time of the request being made and the maximum session start time,  $T_{s_{max}}$  is fixed at 84600s or 24 hours from the time of the request being made. The experiments are simulated over 200,000s or 55.55 hours simulation time. To ensure the simulation is in a stable state,  $t_1$  is set at 100,000s and  $t_2$  is set at 200,000s. Each simulation was run using different random number generator seeds in order to investigate the deviation caused by the simulation model. The data are then used to plot the confidence intervals i.e. mean, maximum and minimum values of the results.

### 3.2 Bandwidth Negotiation At Resource Limit (BNRL)

In this scheme, negotiation only takes place when the requested bandwidth ( $b_r$ ) is not available at the RRT, from  $T_{s_r}$  to  $T_{s_r} + T_r$ . Rather than having a request being rejected, the user may be willing to tolerate a certain degree of service quality degradation by lowering the bandwidth requirement. For example, the user may not mind getting a ‘high quality’ video session (e.g. 10Mbps) if ‘premium quality’ video session (e.g. 12 Mbps) cannot be granted. This tolerance value,  $b_{tol}$  is defined as the percentage of bandwidth reduction tolerable, or  $b_{tol} = (b_r - b_g / b_r) * 100\%$ , where  $b_g$  is the bandwidth granted. The service request utility function with bandwidth negotiation can be represented by  $u(b_r - b_{tol} * b_r, T_{s_r}, T_r, P_r)$ . The offered load  $\% \bar{R}_q$  is also defined as percentage mean load requested in relative to the total bandwidth capacity of the link ( $C$ ). In other words, it represents the mean percentage reservation load at RRT if all the requests have been accepted regardless of the link capacity.

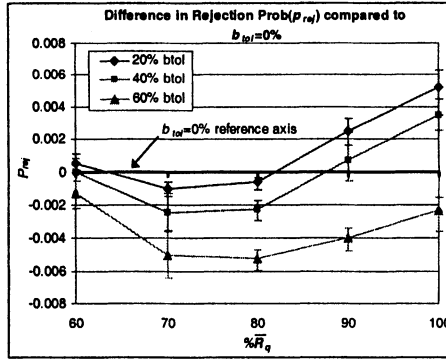


Figure 3. Difference in  $\rho_{rej}$  vs.  $\bar{R}_q$

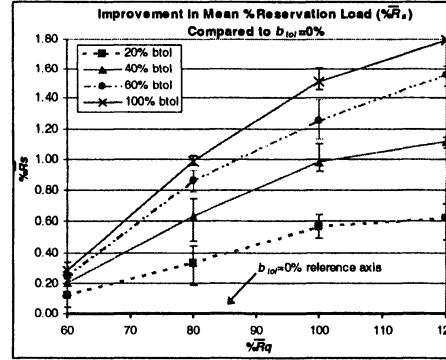


Figure 4.  $\bar{R}_s$  improvement vs.  $\bar{R}_q$

Figure 3 illustrates the difference in  $\rho_{rej}$  when comparing with bandwidth negotiation ( $b_{tol} = 20\%$ ,  $40\%$ , and  $60\%$ ) and without bandwidth negotiation ( $b_{tol}=0\%$ ). In this figure, the x-axis represents the  $b_{tol} = 0\%$  reference point. It can be deduced that at lower  $\bar{R}_q$  (60-80%), bandwidth negotiation generally reduces the  $\rho_{rej}$ . However, when  $\bar{R}_q$  is at 90%,  $b_{tol} = 20\%$  and  $40\%$  experienced higher  $\rho_{rej}$  as compared to  $b_{tol}=0\%$ . The reason is that as  $b_{tol}$  increases, the mean reservation load  $\bar{R}_s$  or RRT utilisation also increases (refer figure 4). Therefore, less bandwidth is available for the next incoming requests and this causes the rejections to increase. As shown in figure 4, the overall  $\bar{R}_s$  improves as  $b_{tol}$  increases. Here, the difference in  $\bar{R}_s$  is almost negligible at low offered load because bandwidth negotiations only happen at high  $\bar{R}_q$ . An improvement of 1  $\bar{R}_s$  means 1 Mbps of extra bandwidth reserved on average. However, this reaches saturation when  $\bar{R}_q$  exceeds 100%. This is due to the higher probability of blocking experienced.

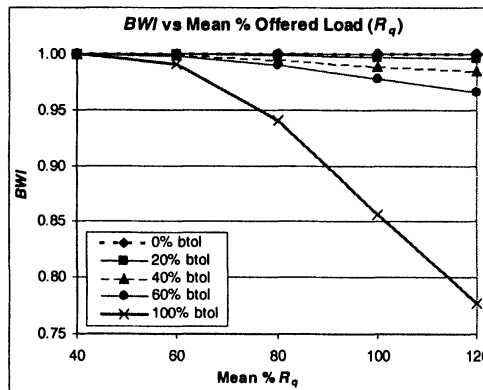


Figure 5. BWI vs.  $\bar{R}_q$

The notion of Bandwidth Index (BWI) is introduced in figure 5. This is defined as the mean ratio of bandwidth granted over bandwidth requested or  $BWI = \text{mean}(b_g/b_r)$  for all the users. It may be viewed as corresponding to the overall users' satisfaction. The figure shows that at low offered load, most users get what they want ( $BWI \sim 1$ ). However at high offered load, the BWI drops as more users have their bandwidth discounted. However, it is noticeable that for  $b_{tol}$  up to 60%, the users still get more than 95% of their original requested bandwidth  $b_r$  (or 0.95 BWI) on average even at 120%  $\bar{R}_q$ . It is important for the network provider to make sure that the overall user satisfaction does not fall below a certain threshold while maximising RRT utilisation and revenue.

Table 1. Improvement in gross revenue ( $Rev_{total}$ ) per hour

$\% \bar{R}_q$	Bandwidth Tolerance ( $b_{tol}$ )			
	20%	40%	60%	100%
60	\$0.24	\$0.40	\$0.49	\$0.57
80	\$0.66	\$1.26	\$1.72	\$1.97
100	\$1.13	\$1.96	\$2.51	\$3.02
120	\$1.24	\$2.22	\$3.12	\$3.59

Table 1 shows the extra gross revenues earned per hour when bandwidth negotiation is possible. This is assumed that 1 Mbps of guaranteed bandwidth is priced at \$2 per hour. Therefore the higher the utilisation, the higher the revenue generated. We can see that even with  $b_{tol} = 20\%$ , extra revenue of \$0.66 per hour can be generated from a 100Mbps link at  $\% \bar{R}_q = 80$  and up to \$3.59 extra revenue per hour can be generated with 100%  $b_{tol}$ .

### 3.3 Bandwidth Negotiation At Predefined Load Level (BNPLL)

An alternative negotiation scheme is proposed where the NPA may initiate negotiation even if the requested bandwidth ( $b_r$ ) is available at RRT. The scheme is set in such a way that when the mean percentage reservation load at RRT during session duration  $T_r$  ( $\% \bar{R}(T_r)$ ) exceeds a predefined load level called Start Negotiate Level (SNL), a bandwidth reduction of  $x\%$  from  $b_r$  will be proposed. However if the user is unwilling/unable to accept the proposed bandwidth ( $>b_{tol}$ ), the original requested bandwidth would be granted anyway.

It can be observed from figure 6 that the BNPLL scheme further reduces the  $p_{rej}$  as compared to BNRL scheme. Essentially, the lower the SNL, the lower the  $p_{rej}$  and the effect is more significant at high  $\% \bar{R}_q$ . This is because the NPA has virtually 'persuaded' all the users to reduce their bandwidth requirements. Hence, more bandwidth is available for future users.



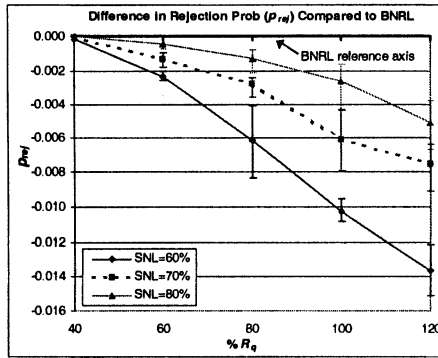


Figure 6. Difference in  $\rho_{rej}$  vs.  $\%R_q$

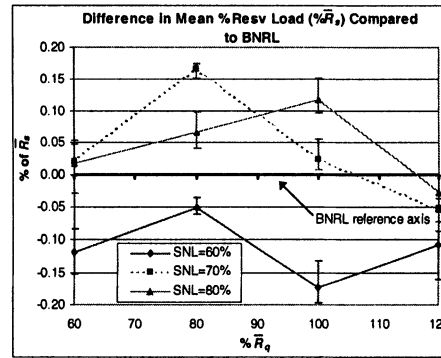


Figure 7.  $\%R_s$  Improvement vs.  $\%R_q$

Figure 7 shows that if SNL is set at 70%, a further improvement of  $\sim 0.17\%R_s$  can be obtained at  $80\%R_q$ . However, if the start negotiate level is set too low (SNL=60%), this will result in lower overall  $\%R_s$ . This shows that higher RRT utilisation can be achieved if the right SNL is applied. Of course, the users must be willing to tolerate bandwidth degradation in the first place.

From figure 8, it is shown that BNPLL results in the significant Bandwidth Index (BWI) degradation although most values still remain above 0.9. Therefore an optimum SNL must be found in order to reach a balance between customers' satisfaction, service availability and RRT utilisation or revenue.

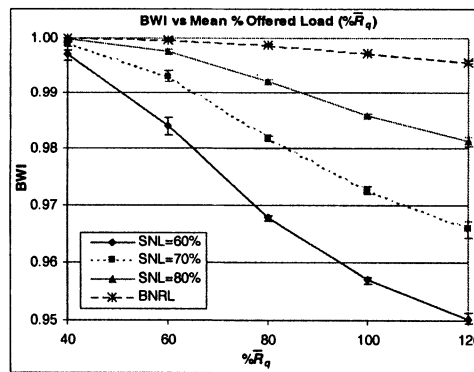


Figure 8. BWI vs.  $\%R_q$

Table 2 presents the difference in gross revenues per hour earned with BNPLL scheme as compared to BNRL scheme if 1 Mbps of guaranteed bandwidth is priced at is \$2 per hour.

Table 2. Improvement in  $Rev_{total}$  per hour compared to BNPL

$\% \bar{R}_q$	Start Negotiate Level (SNL)		
	at 60%	at 70%	at 80%
60	-\$0.24	\$0.04	\$0.03
80	-\$0.10	\$0.33	\$0.13
100	-\$0.34	\$0.05	\$0.23
120	-\$0.21	-\$0.11	-\$0.05

It is shown that if SNL is set at 70%, an extra revenue of \$0.33 per hour can be generated at 80%  $\bar{R}_q$  compared to BNPL scheme with the same offered load and  $b_{tot}$ . Since the drop in revenue (or utilisation) is quite significant with SNL at 60%, it is therefore not advisable to negotiate too early even though this increases the service availability

#### 4. RELATED WORK

The BRPA architecture is in agreement with the policy framework proposed by IETF Policy Framework Working Group [9] and Resource Allocation Protocol Working Group [10], where the immediate goal is to manage QoS provisioning on top of existing QoS technologies. In BRPA architecture, the higher-level agents such as UA, SPA, NA, CPA are essentially the Policy Decision Point (PDP) agents and the element layer router agents (RAs) are the policy enforcers Policy Enforcement Point (PEP) agents.

The concept of resource reservation in advance has also been addressed in [11], [12], [13], [14], etc. To our knowledge, none of these works provides detailed analysis on session-level negotiation involving bandwidth. [11] supports advanced reservation based on predictive service admission control algorithm. It also presents some initial considerations concerning the mapping of advanced reservation protocol onto RSVP. Work by [13] focuses on the design, implementation and evaluation of their Resource Reservation in Advance (ReRA) mechanism by extending the existing RSVP protocol on ATM. The authors also address best-match alternative reservation scenarios similar to that offered by BRPA bandwidth negotiation schemes. However in their work, no detailed negotiation algorithms and experiments are provided. [12] presents a general architecture that describes the requirements of ReRA. A simple prototype that employs some user agents has been developed. The basic tasks of their user agents are very similar to BRPA' UA's that includes performing reservation, control or modify the state of existing reservations, acknowledgement, announcement, etc. [14] proposes an agent-based reservation system for immediate and advance reservation calls. In their work, a call 'lookahead' time is applied to decide the admission of IR calls.

The effects on rejection probability, pre-emption probability and overall RRT utilisation are studied.

## 5. CONCLUSIONS

The adaptive and distributive nature of the agent enhances greater flexibilities and automation in service provisioning, especially in today's distributed, heterogeneous, and fast changing network environment. The main objective of this chapter is to devise a policy-enabled, flexible resource provisioning architecture that facilitates SLA negotiations and brokering. The Bandwidth Resource Provisioning Architecture (BRPA) architecture has been introduced to facilitate distinct end-to-end path with bandwidth, session duration, session start time preferences negotiations, etc on a per-user, per-application or per-flow basis. In this paper however, due to the limited space only session-level negotiations schemes involving bandwidth have been presented. Further analyses on session start time, duration negotiations are provided in [15].

The results show that these schemes can be exploited for the benefits of both negotiating parties. It is shown that in most cases, negotiation reduces rejection probability and improves mean RRT utilisation and therefore network's revenues. Choosing which scheme to be applied depends very much on the type of applications, the user's preferences and also the load of the link during the time of negotiation. Some strategies need to be devised when implementing these schemes so that overall optimisation can be achieved. For example, different service qualities can be offered to dictate the mean bandwidth request. Various policies can also be applied in the future to control the session durations, session start time boundary ( $T_{s_{min}}$  and  $T_{s_{max}}$ ) etc. Some pricing strategies can also be applied to control the users' behaviours. Indirectly, these are seen as a means to manage bandwidth resource.

## References

- [1] G. Huston, "Next Steps for the IP QoS Architecture", *IETF Internet Draft (draft-iab-qos-02.txt)*, Aug 2000.
- [2] David Chieng, Alan Marshall, Ivan Ho, Gerald Parr, "A Mobile Agent Brokering Environment for The Future Open Network Marketplace", *Seventh International Conference On Intelligence in Services and Networks (IS&N2000)*, pp. 3-15, Athens, 23-25 February 2000. (Springer Verlag LNCS Vol. 1774)
- [3] Q. Ma and P. Steenkiste, "Quality of Service Routing for Traffic with Performance Guarantees", *IFIP International Workshop on Quality of Service (IWQoS'97)*, New York, May 1997, pp. 115-126.

- [4] R. Braden, L. Zhang, S. Berson, S. Herzog and S. Jamin. "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification", *RFC 2205*, IETF Network Working Group, Sept 1997.
- [5] C. Topolcic, "Experimental Internet Stream Protocol, Version 2 (ST-II)", *RFC 1190*, IETF, Oct 1990.
- [6] P. Srisuresh and P. Joseph, "TE LSAs to extend OSPF for Traffic Engineering", *IETF Internet Draft (draft-srisuresh-ospf-te-02.txt)*, Jan 2002.
- [7] L. Mathy, C. Edwards, and D. Hutchison, "The Internet, A Global Telecommunications Solution?", *IEEE Network*, Vol. 14, No. 5, pp. 46-57, July/Aug 2000.
- [8] David Chieng, Alan Marshall, Ivan Ho and Gerald Parr, "Agent-Enhanced Dynamic Service Level Agreement In Future Network Environments", *IFIP/IEEE MMNS 2001*, pp. 299-312, Chicago, 29 Oct - 1 Nov 2001. (Springer Verlag LNCS Vol. 2216)
- [9] IETF Policy Framework Working Group. ([www.ietf.org/html.charters/policy-charter.html](http://www.ietf.org/html.charters/policy-charter.html))
- [10] The Resource Allocation Protocol Working Group.
- [11] Mikael Degermark et. al., "Advance Reservation for Predictive Service in the Internet", *ACM/Springer Verlag Journal on Multimedia Systems*, Vol. 5, No. 3, pp. 177-186, May 1997.
- [12] Lars C. Wolf and Ralf Steinmetz, "Concepts for Resource Reservation in Advance", *Special Issue of the Journal of Multimedia Tools and Applications on "The State of The Art in Multimedia"*, Vol. 4, No. 3, May 1997.
- [13] Alexander Schill, Frank Breiter and Sabine Kuhn, "Design and Evaluation of an Advance Reservation Protocol on top of RSVP", *IFIP 4th International Conference on Broadband Communications*, pp. 23-24, Stuttgart, March 1998.
- [14] O. Schelen and S. Pink, "Resource sharing in advance reservation agents", *Journal of High Speed Networks: Special issue on Multimedia Networking*, Vol. 7, No. 3-4, pp. 213-28, 1998.
- [15] David Chieng, "A Framework for Provisioning Network Resources based on Agent-Enhanced Service Level Agreements", *PhD Thesis*, School of Electrical and Electronic Engineering, The Queen's University of Belfast, May 2002.