# **Evaluation of a Scalable Reservation-Based QoS Architecture**

Rui Prior<sup>1</sup>, Susana Sargento<sup>1,2</sup>, Pedro Brandão<sup>1</sup>, Sérgio Crisóstomo<sup>1</sup> <sup>1</sup>DCC/FCUP & LIACC, University of Porto <sup>2</sup>Institute of Telecommunications, University of Aveiro {rprior, ssargento, pbrandao, slc}@ncc.up.pt

## Abstract

This paper presents a performance evaluation study of the Scalable Reservation-Based QoS architecture. This architecture introduces a scalable per-flow signalling model, using techniques and algorithms developed to minimise the computational complexity, namely a label switching mechanism and an efficient timer implementation. The underlying architecture is based on DiffServ and resource reservation is performed for aggregates of flows at both core and access networks. The obtained results show that this architecture is able to provide QoS guarantees, irrespectively of the behavior of other flows in the same and in different classes, maximizing the network resource utilization. Based on the performance evaluation, we can state that this architecture is able to support service classes with strict and soft QoS guarantees in high speed networks.

## 1. Introduction

Nowadays, the Internet does not support service models other than best effort, essential for the transport of emerging applications. With the objective of providing the Internet with Quality of Service (QoS) and differentiation, the IETF proposed two main QoS architectures. The Integrated Services (IntServ) architecture [1] provides strict QoS guarantees and efficient resource usage, but suffers from scalability problems, concerning the per-flow scheduling, classification and reservation procedures. The Differentiated Services (DiffServ) architecture [2], based on aggregation, is free from these scalability concerns, but without an admission control mechanisms to limit the number of flows in the network, all flows belonging to a class may be degraded.

Aiming at the introduction of QoS support without the aforementioned limitations, several architectures have been proposed in the literature [3, 4, 5, 6, 7, 8]. All of these architectures, however, suffer from one or more of the following problems: lack of strict QoS guarantees, underutilization of network resources, or scalability limitations steming from the complexity of the algorithms and procedures

used, among others. In [9] we proposed an architecture, denoted Scalable Reservation-Based QoS (SRBQ), that provides end-to-end QoS support without the problems of the previously mentioned architectures. It does not impose a complex scheduling mechanism, uses efficient aggregatebased packet classification, supports classes with both soft and strict QoS guarantees and achieves a good utilization of network resources.

In this paper we present a performance evaluation study of the SRBQ architecture, based on simulation results. We analyze the standard QoS parameters (delay, jitter and packet losses) and the network utilization in different experiments, using both synthetic flows and real-world multimedia streams. In terms of QoS guarantees, this paper shows that the service classes in our architecture provide strict and soft QoS guarantees, irrespectively of the behavior of other flows in the same and in different classes. The behavior of the architecture with real traffic flows reinforces the QoS guarantees achieved with SRBQ. In terms of resource utilization, we also show that all these guarantees are achieved with high network utilization.

This rest of the paper is organized as follows. Section 2 contains a brief overview of the SRBQ architecture. In section 3 the performance results of SRBQ are analysed. Finally, section 4 contains the most important conclusions, and describes the future work to be performed and extensions to be applied to the architecture.

#### 2. Architecture Overview

The underlying architecture of SRBQ is based on Diff-Serv, with the addition of signaling-based reservations subject to admission control. The network is partitioned into domains, consisting of core and edge routers; access domains contain also access routers. Flows are aggregated according to service classes, mapped to DiffServ PHBs (Per-Hop Behaviors), and packet classification and scheduling are based on the DS field of the packet headers. Besides Best-Effort (BE), SRBQ supports a Guaranteed Service (GS) class, providing strict QoS guarantees, and one or more Controlled Load (CL) classes, emulating lightlyloaded BE networks, based on the Assured Forwarding (AF) PHB.

SRBQ's queuing model is compatible with DiffServ; the two models may coexist in the same network. The main scheduler is priority-based: the highest priority belongs to the GS class, which must be shaped by a token-bucket; below, there is a class for signaling, which must be ratecontrolled; the CL class(es) come next, with optional ratecontrol; finally, at the bottom priority is the BE class.

In SRBQ, all nodes perform signaling and support the previously described queuing model. Access routers perform per-flow policing for the CL class and per-flow ingress shaping for the GS class. Edge routers perform aggregate policing and DSCP remarking. Core routers perform no policing.

Flows are subject to admission control, performed at every node. GS flows are characterized by token-buckets; CL flows are characterized by 3 rate water-marks, corresponding to different drop priorities. A scalable hop-byhop signaling protocol was developed to perform unidirectional, sender initiated, soft-state reservations. It uses a label switching mechanism, developed to allow direct access to the reservation structures, and an efficient implementation of soft timers. The labels are installed at reservation setup time, and all subsequent signaling messages use them.

## 3. Performance Results

The SRBQ architecture has been implemented in the ns-2 simulator. In this section we present the performance results obtained by simulation, which mainly address the QoS guarantees achieved by the model. Though very important in SRBQ, processing efficiency measurement is out of the scope of this paper, since the ns-2 simulator is not suited for the evaluation of processing delays. The considerations on the scalability of the model presented in [10], though, indicate that it is suitable for use in high-speed core networks.

The simulated scenario is depicted in figure 1. It includes 1 transit and 5 access domains. Each terminal in the access domains simulates a set of terminals. The reason for having more than one access domain connected to an edge node of the access and transit domains is to check that correct aggregate policing is performed at the entry of the domain. The bandwidth of the connections in the transit domain, and in the interconnections between the transit and the access domains, is 10 Mbps. The propagation delay is 2 ms in the transit domain connections and 1 ms in the interconnections between the access and the transit domain.

In this scenario we consider the coexistence of GS, CL and BE classes. At each referred connection, the bandwidth assigned to the signalling traffic is 1 Mbps. Note that, although this seems very high, the unused signalling bandwidth is used for BE traffic. Except where otherwise stated, the bandwidth assigned to the GS class is 3 Mbps, while for CL it is 4 Mbps. The remaining bandwidth is used for BE traffic. The bandwidth reserved for the GS and CL classes and left unused is also used for BE.



Figure 1. Simulated scenario

Each terminal of the access domains on the left side generates a set of flows belonging to the GS, CL and BE classes. Each source may generate traffic to all destinations; the destination of each flow is randomly chosen in the set of the terminals in the right side access domains. Traffic in each class is a mixture of different types of flows.

Both parameter-based (PBAC) and measurement-based (MBAC) admission control are supported in the CL class. The same sets of simulations were performed using both methods. The algorithm used for MBAC is a modification of Measured Sum (MS) with support for 3 different levels, corresponding to the water-marks, and a global target utilization of 95% was used.

All simulations presented in this paper are run for 5,400 simulation seconds, and data for the first 1,800 seconds is discarded. All values presented are an average of, at least, 5 simulation runs with different random seeds. The next sub-sections present the results of these experiments.

#### 3.1. End-to-End QoS Guarantees

In this sub-section we discuss the results of 3 sets of experiments used to evaluate the QoS performance of the SRBQ architecture. The set of flows is distributed in the following way (table 1): (1) traffic in the GS class is composed by CBR (Constant Bit Rate) flows (Voice and Video256) and on-off exponential (Exp1gs) flows; (2) traffic in the CL class is composed by on-off exponential (Exp1cl) and Pareto (Pareto1cl) flows; and (3) traffic in the BE class is composed by on-off Pareto (Pareto1be) and FTP (Ftpbe) flows. Flows belonging to the BE class are active for the overall duration of the simulations (there are 3 FTP and 2 Pareto flows per source), while flows in the other classes are initiated according to a Poisson process with a certain mean time interval between calls (MTBC), and each flow has an average duration (Avg dur.) exponentially distributed.

The largest Mean Offered Load (MOL) in the GS and CL classes is, in terms of average traffic rates, about 20% higher than the bandwidth assigned to those classes, which, due to different mixes of flow types, translates, in terms of requested reserved rates (ROL - Requested Offered Load), in excess figures of 26% (GS) and 42% (CL). The values presented in the table correspond to this maximum offered load, denoted by a load factor of 1. For lower amounts of offered traffic, the MTBC is increased in the inverse proportion of the offered load factor.

For GS flows, the reservation rate (Resv rate) represents the rate of the token bucket and the reservation burst (Resv burst) represents its depth. The reservation parameters provide a small amount of slack to compensate for numerical errors in floating point calculations. For CL flows, Low RR (Reservation Rate), Resv rate and High RR represent the three rate water-marks used for drop precedence selection and packet dropping at the policer. The same sets of simulations were performed using MBAC and PBAC in the CL class. The utilization limits for the three rate water-marks were set to 0.7, 1.0 and 1.7 times the bandwidth assigned to this class. The sum of the rates in each water-mark for all flows in the class must not exceed the respective utilization limits. Notice that both scheduling and policing are performed on a per-class basis (except at the access routers).

In the first set of experiments we evaluate the end-toend QoS guarantees of both GS and CL classes for different amounts of CL offered traffic. Figures 2 (a and b) present the packet loss ratio and mean utilization of both GS and CL flows when the offered load factor of the GS flows is 1 and the offered load factor of the CL flows increases from 0.5 to 1, combining the results of simulations with PBAC and MBAC in the CL class. Notice that GS always uses PBAC. Packet losses are null for well behaved GS flows. In CL flows, packet losses increase with the offered load, but remain nevertheless very low (less than 0.03%) when using PBAC. This means we could be more aggressive by reducing the requested rate water-marks for these flows. With MBAC, losses raise significantly, due to the higher utilization figures, reaching 0.25% for the heavy-tailed Pareto flows. Losses for exponential GS flows, of about 0.13%, are due to buffer space limitation at the ingress shaper, since these flows are not conformant to their reservations. At the core, the average utilization of the GS class is just below 2.5 Mbps (83%), and that of the CL class varies, with a decreasing slope, from 2.4 Mbps (60%) to 3.1 Mbps (78%) with PBAC, and from 2.4 Mbps (60%) to 3.3 Mbps (83%) with MBAC.

The average delay (not shown) remains very low and almost constant for all flow types, except for the GS exponential flows. For all except these, the delay is mostly the sum of transmission and propagation delays. GS exponential flows suffer an additional, and potentially large, delay at the ingress shaper of the access router when they send at a rate larger than what they requested for long periods of time. It is the applications' fault, though, for transmitting non-conformant traffic. The fact that the delay for the other GS flows remains very low shows that they are not adversely affected. The delay for CL flows remains almost constant, independently of the offered traffic. Jitter figures exhibit a similar behavior as delay one.



Figure 2. Loss and utilization with varying CL offered load

In the second set of experiments we varied the offered load of the GS flows from 0.5 to 1, keeping a constant CL load factor of 1. The QoS results in both GS and CL flows are not affected by the GS offered load.

In the third set of experiments we analyse the effect, on the delay and packet losses of both GS and CL classes, of decreasing the requested rate of the CL flows. Figures 3 (a and b) show, respectively, the variation of the delay and packet loss values with varying requested rates for CL flows. Here we have set the flow acceptance utilization limits of the three rate water-marks to 0.7, 1.0 and 2.0 times the bandwidth assigned to CL in order to ensure that flow admission would be performed based on the second rate water-mark, the varying factor in these experiments. Since the average rate for both types of CL flows is 128 kbps, we vary the requested rate from 130 kbps to 160 kbps, a little higher than the 150 kbps used in the previous experiments. As a result of increasing the reserved rate, the average uti-

Class	Туре	Peak rate	On time	Off time	Avg. rate	Pkt size	Resv rate	Resv burst	Low RR	High RR	MTBC	Avg dur.	MOL	ROL
		(kbps)	(ms)	(ms)	(kbps)	(bytes)	(kbps)	(bytes)	(kbps)	(kbps)	(s)	(s)	(kbps)	(kbps)
	Voice	48	-	-	48	80	48.048	81	-	-	45	120	768	769
GS	Video256	256	-	-	256	1000	256.256	1050	-	-	180	240	2048	2050
	Exp1gs	256	200	200	128	1000	160	5000	-	-	90	90	768	960
<u></u>	Pareto1cl	256	200	200	128	1000	150	-	64	256	38	120	2425	2842
UL.	Exp1cl	256	200	200	128	1000	150	-	64	256	38	120	2425	2842
											Simu	t. Flows		
DE	Ftpbe	-	-	-	-	1040	-	-	-	-	3 per sr	c terminal	var.	N.A.
DE	Paroto1bo	256	200	200	128	1000	-		-	-	2 nor cr	c termina	2304	ΝΔ

Table 1. Characteristics of the traffic flows (end-to-end QoS tests)

lization of the CL class at the core decreases from 3.5 Mbps (88%) to 3.0 Mbps (75%) with PBAC, and from 3.5 Mbps (88%) to 3.2 Mbps (80%) with MBAC. As expected, the difference in utilization of the CL class with PBAC and MBAC is higher when the difference between the reserved rate and the actual transmission rate is larger.

GS flows are not affected by the CL traffic: the delay for CBR GS flows remains constant, and is approximately equal to the sum of transmission and propagation delays; exponential GS flows experience a much higher delay due to the ingress shaper. As expected, the delay for CL flows decreases with the increasing requested rate, since the number of accepted flows is lower, and less packets are being marked for the highest drop probability. Jitter figures, though not shown, have a similar variation pattern. The most interesting results for this group are the loss figures. Packet loss in GS flows is not affected by the CL reservations, being null for conformant flows. CL flows, on the other hand, exhibit increasing losses with decreasing requested rates. With PBAC and a requested rate of 130 kbps, which is only 1.6% higher than the average transmission rate, packet loss for exponential CL flows is just below 0.5%, while for the heavier tailed Pareto it is slightly above 0.8%. With MBAC, these values raise to 0.9% and 1.3%, respectively. With reservations of 160 kbps, CL losses are below 0.1% with MBAC and below 0.007% with PBAC, even for the heavy-tailed Pareto flows.

These sets of experiments show that our model, though being aggregation-based, is able to support both strict and soft QoS guarantees and achieves complete independence between traffic classes.

#### 3.2. Independence Between Flows

In this sub-section we evaluate the performance of the architecture in the presence of misbehaved flows, that is, flows that send at a rate much higher than the one they requested for considerable periods of time. Moreover, we also analyse the influence of misbehaved flows on well behaved ones in both the GS and the CL classes. In order to protect the network from non-conformant flows, the access router performs per-flow ingress shaping for GS class flows. This shaper absorbs multiplexing jitter from the terminal and ensures that the traffic injected into the network does not exceed the reserved parameters by absorbing ap-



a) Mean delay vs increasing reserved rate



Figure 3. Delay and packet losses with varying reserved rates for the CL flows

plication bursts above the requested bucket (of 5 packets in this case), thus protecting the other GS flows. CL flows, on the other hand, are policed, instead of shaped, at the ingress router. This means that a single misbehaved CL flow will be penalized in terms of packet losses, but will not be significantly affected in terms of delay.

In this experiment, the mean offered load (MOL) for both classes is 23% larger than their assigned bandwidth (table 2). There are three types of flows in each class: (1) a CBR flow (Video64) that is considered a well behaved flow; (2) an on-off exponential flow (Exp1) with a burstiness of 50% (average busy and idle times of 200 ms) and a peak rate of 256 kbps, that is considered a nearly well behaved flow, since it sends at a rate a little higher than requested; and (3) an on-off exponential flow (Exp2) with varying burstiness and peak rate, that is considered a misbehaved flow, since it sends at a rate much larger than requested for considerable periods of time. Its burstiness is variable, from 50% to 12.5%, varying its peak rate between 256 kbps (average

	Class	Type	Peak rate	On time	Off time	Avg. rate	Pkt size	Resv rate	Resv burst	Low RR	High RR	MTBC	Avg dur.	MOL	ROL	
			(kbps)	(ms)	(ms)	(kbps)	(bytes)	(kbps)	(bytes)	(kbps)	(kbps)	(s)	(s)	(kbps)	(kbps)	
ſ		Video64gs	64	-	-	64	500	64.064	501	-	-	75	240	1229	1230	L
	GS	Exp1gs	256	200	200	128	1000	160	5000	-	-	75	120	1229	1536	Ĺ
		Exp2gs	var.	var.	var.	128	1000	160	5000	-	-	75	120	1229	1536	
ſ		Video64cl	64	-	-	64	500	65	-	64	66	75	240	1229	1248	L
	CL	Exp1cl	256	200	200	128	1000	150	-	64	256	50	120	1843	2160	
		Exp2cl	var	var	var	128	1000	150	-	64	256	50	120	1843	2160	í.

Table 2. Characteristics of the traffic flows (isolation tests)

busy and idle times of 200 ms) and 1024 kbps (average busy and idle times of 50 ms and 350 ms, respectively). Notice that the sum of the average idle and busy times remains constant (400 ms), as does the average rate. It is the high mismatch between the requested rate and the peak transmission rate that turns Exp2 flows into misbehaved ones. In this test, only PBAC is used.

Figures 4 (a and b) depict the packet loss ratio and the mean delay for all three types of flows with increasing burstiness values of the misbehaved (Exp2) flows. There are no packet losses for well behaved (Video64) GS flows; losses for nearly well behaved (Exp1) GS flows are just above 0.1%. Packet losses for misbehaved (Exp2) GS flows reach 7.1% when their burstiness reaches 12.5%. With such a burstiness, the peak rate of this type of flow is much larger than the reserved rate, and a large number of packets is lost. Their misbehavior, however, does not affect the previous flows.

The CL class does not provide the same absolute guarantees as the GS class: though not easily seen in the figure, losses vary from 0.015% to 0.002% for well behaved (Video64), and from 0.024% to 0.006% for nearly well behaved (Exp1) CL flows, when the burstiness of the misbehaved (Exp2) CL flows varies from 12.5% to 50%. We realize that even nearly well behaved flows experience very small losses. Losses in misbehaved flows vary from 4.6% to 0.006%; they are penalized for their burstiness.

The mean delay of the well behaved GS flows is very small, and is mainly due to transmission and propagation delays. Nearly well behaved GS flows have a constant average delay in the order of 160 ms, which is significantly larger than that of the well behaved ones. Notice that this type of flow has a peak bandwidth approximately 100 kbps larger than the requested one, and therefore the packets will experience some delay (and small amounts of losses) at the ingress shaper of the access routers when the sources transmit at the peak rate for longer periods of time. As expected, misbehaved GS flows have a delay that increases with their burstiness: with a burstiness of 12.5%, this delay can reach more than 400 ms. GS jitter curves, though not shown, exhibit the same behavior as their delay counterparts. Notice that since all GS flows are aggregated and use the same queue, internally served in a FIFO fashion, the queueing delay is shared by all GS flows. Therefore the large delays for nearly well behaved and misbehaved GS flows are inflicted at the ingress shaper. This reaction against misbehaved

flows (in terms of large delays and losses) is meant to protect the other GS flows. This way, well behaved GS flows preserve a constant and small delay and no packet losses irrespectively of the burstiness of the misbehaved flows. It is the applications' fault for requesting inadequate reservations in face of the traffic to be transmitted.

Regarding delay, the behaviour of the CL class is entirely different from that of the GS class. Since there is no shaping at the access router, only policing, no significant delay penalty is inflicted to misbehaved flows: on average, it is only 2.6% higher with a burstiness is 12.5% than it is with a burstiness of 50%. There is also, however, a very slight increase of about 1.5% in the average delay of other flows belonging to this class in the presence of misbehaved flows.





This experiment shows that the system reacts accordingly in the presence of misbehaved flows, keeping a complete independence between GS flows, which is expected in this type of service. The CL class is more tolerant to bursty flows: since there is no shaping, they are not penalized in terms of delay, only in terms of packet losses inflicted by the access and edge routers by means of packet (re)marking and policing. These losses ensure that network congestion remains low, protecting (though not completely isolating) conformant CL flows. The CL class is, therefore, more appropriate for misbehaved flows with soft QoS requirements.

These QoS results are not unusual: the main achievement of our model is to provide them in a scalable, aggregationbased architecture, while keeping a good utilization of network resources.

## 3.3. Real multimedia streams

All the previous experiments were based on synthetic flows with specific characteristics meant to evaluate particular aspects of the performance of the SRBQ architecture. In order to evaluate its performance under normal working conditions we performed a set of tests using packet traces from real multimedia flows. The trace files we used correspond to H.263 video streams with average bit rates ranging from 16 kbps to 256 kbps, and are available from [11].

In these experiments, we have assigned 2 Mbps to the GS class and 5 Mbps to the CL class. The rate water-marks for the CL class were adjusted to 0.8, 1.0 and 2.0 times the bandwidth assigned to it. Table 3 summarizes the parameters of the flows. Traces from several different video streams were used for each bit rate, and the starting point in the stream for each flow was randomly chosen. Flows are initiated according to a Poisson process and have a duration following a Pareto distribution. The highest mean offered load for both the GS and the CL classes is 20% higher than their assigned bandwidth, which is denoted by a load factor of 1. These simulations were performed using both PBAC and MBAC in the CL class.

Figure 5 shows the performance and utilization results for varying offered load factors in both classes, combining results from the simulations using PBAC and MBAC in the CL class; PBAC is always used in the GS class. Delay in the GS class does not seem to be affected by the offered load factor, while that of the CL class exhibits a slight growth trend, more evident when using PBAC. Notice that the delay is always smaller than 20 ms. Jitter figures have a similar behaviour. The higher jitter values in the GS class are due to ingress shaping at the access router, performed in order to force the flows into conformance with the reservations.

There are no packet losses in the GS class: burstiness above the reserved rate is absorbed by the ingress shaper at the access router (within certain bounds), and translates into increased delay and jitter rather than losses. Packet loss curves for the CL class exhibit a seemingly contradictory behaviour: they are higher for lower values of offered load. There is, however, an explanation for this fact, which stems from a combination of several factors. (1) Lower bit rate h.263 flows are more bursty and have smaller packets



Figure 5. Performance and utilization results using real video flows

than higher bit rate ones. (2) The second rate water-mark (for which the class has a target utilization factor of 1) used in the reservations for these flows is 3% higher than their target bit rate. Therefore, the absolute difference between the reserved rate and the target bit rate is larger in higher

Class	Туре	Avg rate	Pkt size	Resv rate	Resv burst	Low RR	High RR	MTBC	Avg dur.	MOL	ROL
		(kbps)	(bytes)	(kbps)	(bytes)	(kbps)	(kbps)	(s)	(S)	(kbps)	(kbps)
	video	16	var.	17	4000	-	-	151	180	114	122
GS	video	64	var.	68	5000	-	-	151	180	458	486
	video	256	var.	272	8000	-	-	151	180	1831	1945
	video	16	var.	16.5	-	12	40	60	180	288	297
CL	video	64	var.	66	-	48	96	60	180	1152	1188
	video	256	var.	264	-	192	352	60	180	4608	4752
								Simult	Simultan. flows		
DE	ftp	-	-	-	-	-	-	3 per si	rc terminal	var.	N.A.
DE	noroto	100	1000					2 por o	torminal	1526	NI A

#### Table 3. Characteristics of real traffic flows

bit rate flows. The absolute difference between the third water-mark and the target bit rate is also larger in higher bit rate flows. (3) Packet remarking and policing is performed on an aggregate basis at the edge routers. These 3 factors combined mean that when there are more high bit rate flows in the network, low bit rate ones take advantage of the bandwidth excess from the higher bit rate ones, therefore decreasing their loss ratio. This fact has a much higher weight in the overall packet loss ratio than the (minimal) amount of network congestion.

The utilization of each class, as expected, grows with the offered load. With a load factor of 1, the mean utilization of the GS class is 1.3 Mbps (65%), and that of the CL class is 3.8 Mbps (75%) when using PBAC and 3.6 Mbps (72%) when using MBAC. The lower utilization with MBAC than with PBAC has a simple explanation. The reserved rate (second water-mark) for these flows is very close to their target bit rate (only 3% higher). The MBAC utilization factor is 95%, and the algorithm is Measured Sum (MS), which computes the average rate during T time intervals of duration  $\tau$ , using the largest of these values as an estimation of the used bandwidth. This means that most of the time the estimated bandwidth will be higher than 95% of the sum of the reservations, leading to lower utilization figures with MBAC.

The results presented in the previous paragraphs show that the SRBQ architecture is able to meet the QoS requirements of the supported service classes when using realworld multimedia flows.

## 4. Conclusions and Future Work

In this paper we presented a study on the performance evaluation of our QoS architecture, Scalable Reservation-Based QoS (SRBQ). In this study we analyzed the standard QoS parameters and the network utilization in different experiments, using both synthetic flows and real-world multimedia streams. In terms of QoS guarantees, this paper has shown that this architecture is able to provide strict and soft QoS guarantees, irrespectively of the behavior of other flows. These guarantees are achieved with aggregate-based reservations and high network resource utilization. Therefore, this architecture is able to provide both IntServ service models with an underlying DiffServ network, minimizing the processing load at each network element, and yielding good network resource utilization figures.

In order to evaluate and quantify the scalability of the solution as compared to others, we plan to develop a prototype implementation. Further research on the SRBQ architecture will be focused towards the support for accounting and charging, and interoperation with QoS routing protocols.

## References

- R. Braden, D. Clarck, and S. Shenker, "Integrated Services in the Internet Architecture: an Overview," RFC 1633, IETF, June 1994.
- [2] S. Blake, D. Blake, M. Carlson, E. Davies, Z. Wang, and W. Weiss, "An Architecture for Differentiated Services," RFC 2475, IETF, December 1998.
- [3] I. Stoica, Stateless Core: A Scalable Approach for Quality of Service in the Internet. PhD thesis, Carnegie Mellon University, December 2000.
- [4] C. Cetinkaya and E. Knightly, "Egress Admission Control," in *Proceedings of IEEE INFOCOM 2000*, March 2000.
- [5] L. Breslau, E. Knightly, S. Shenker, I. Stoica, and H. Zhang, "Endpoint admission control: Architectural issues and performance," in *Proceedings of ACM SIGCOMM 2000*, August 2000.
- [6] S. Sargento, R. Valadas, and E. Knightly, "Resource Stealing in Endpoint Controlled Multi-class Networks," in *Proceedings of IWDC 2001*, September 2001. Invited paper.
- [7] Y. Bernet, P. Ford, R. Yavatkar, F. Baker, L. Zhang, M. Speer, R. Braden, B. Davie, J. Wroclawski, and E. Felstaine, "A Framework for Integrated Services Operation over Diffserv Networks," RFC 2998, IETF, November 2000.
- [8] F. Baker, C. Iturralde, F. L. Faucheur, and B. Davie, "Aggregation of RSVP for IPv4 and IPv6 Reservations," RFC 3175, IETF, September 2001.
- [9] R. Prior, S. Sargento, S. Crisóstomo, and P. Brandão, "Endto-end Quality of Service with Scalable Reservations," in *Proceedings of ICTSM11*, October 2003.
- [10] R. Prior, S. Sargento, S. Crisóstomo, and P. Brandão, "Efficient Reservation-Based QoS Architecture," LNCS 2899, pp. 168–181, Springer-Verlag, November 2003.
- [11] Arizona State University, "MPEG-4 and H.263 Video Traces for Network Performance Evaluation." http://trace.eas.asu.edu/TRACE/trace.html, 2004.