

Stackelberg approach for pricing Differentiated Services

Eitan Altman* Richard Marquez† Rachid El-Azouzi‡
David Ros§ Bruno Tuffin¶

March 28, 2007

Abstract

We consider in this paper both real-time traffic as well as data transfers sharing a common bottleneck link. We assume that data transfer uses TCP congestion control protocol and that real-time traffic uses some TCP-friendly transport protocol that satisfies the same square-root formula for throughput. The performance measures are determined according to the operational parameters of a RED buffer management. The latter is assumed to be able to give differentiated services to the applications according to their choice of service class. In terms of loss probabilities and of throughputs, we consider a best effort type of service differentiation where the QoS of connections is not guaranteed, but by choosing a better (more expensive) service class, the QoS parameters of a session can improve (as long as that session is the only one to change its service class). We assume however, that the system is dimensioned so as to satisfy some average delay requirement. The choice of a service class of an application will depend both on the utility as well as on the cost it has to pay. We first study the performance of each traffic source as a function of the connections' parameters and the pricing policy of the network. We then study the Stackelberg equilibrium, i.e. the service provider's problem of how to choose the pricing so as to maximize its utility, taking into account the reaction of the users.

Keywords: TCP, Buffer Management, RED/AQM, Stackelberg equilibrium, Pricing

*Address: INRIA, B.P. 93, 2004 Route des Lucioles, 06902, Sophia-Antipolis Cedex.

†Dpto. Sistemas de Control, Univ. Los Andes, Mérida 5101, Venezuela, email: marquez@ula.ve

‡Université d'Avignon et des Pays de Vaucluse (IUP), LIA-CERI, 339 chemin des Meinajariès, B.P.1228, 84911 Avignon Cedex 9, France

§GET/ENST Bretagne, Rue de la châtaigneraie CS 17607, 35567 Cesson Sévigné Cedex, France

¶IRISA/INRIA, Campus Universitaire de Beaulieu, 35042 Rennes Cedex, France

1 Introduction

We study in this paper the performance of competing connections that share a bottleneck link. Both data connections as well as real-time connections are assumed to have controlled rates that are "TCP friendly", i.e. they satisfy the well known square-root relation between marking probability and throughput [10, 12, 13]. A RED buffer management is used for early drop of packets. We allow for service differentiation between the connections through the packet dropping (or marking) probability which may depend on the connection (or on the connection class). More specifically, we consider a buffer management scheme that uses a single averaged queue length to determine the dropping probabilities (similar to the way it is done in the RIO-C (coupled RIO) buffer management, see [14]); for any given averaged queue size, packets belonging to connections with higher priority have smaller probability of being dropped than those belonging to lower priority classes. We compute the throughput and the average drop probability for each connection. We assume in our model that the average queue size is a predetermined system parameter, i.e. there is a target value for the average queue length which implies also a fixed average queueing delay at the buffer. Since the throughputs of connections can vary (according to the price they pay), this means that the link rate is adapted to the connections' input rates so as to guarantee this desirable average queue length.

We then address the question of the choice of priorities. Given utilities that depend on the performance measures, on one hand, and on the cost for a given priority (i.e. the pricing strategy of the provider), on the other hand, each user is faced with an optimization problem, which we solve explicitly. This then allows us to determine the choice of pricing strategy by the service provider which maximizes its own profits; the solution to this bi-level optimization problem is known as the Stackelberg equilibrium.

In a previous paper [3], we analyzed a related problem with TCP and CBR traffic sharing a common bottleneck buffer without a given target on the average queueing delay. This created a strong coupling between the different users, so that the utility of a given user was influenced by the priority choices of all other users. This gave rise to a more complicated modeling of the problem faced by the users (for a fixed pricing strategy of the provider) which was shown to be a non-cooperative game. Although some properties of the equilibrium of the game were obtained in [3], we were not able to obtain explicit formulae for the equilibrium priority choices of the users and therefore did not treat the Stackelberg problem there.

Related references. We briefly mention other recent work in that area. Reference [6] has considered a related problem where the traffic generated by each session was modeled as a Poisson process, and the service time was exponentially distributed. The decision variables were the input rates and the performance measure was the goodput (output rates). The paper restricted itself to symmetric users and symmetric equilibria and the pricing issue was not considered. In this framework, with a common RED buffer, it was shown that an equilibrium does not exist. An equilibrium was obtained and characterized

for an alternative buffer management that was proposed, called VLRED. We note that in contrast to [6], since we also include in the utility of CBR traffic a penalty for losses (which is supported by studies of voice quality in packet-based telephony [8]), we do obtain an equilibrium when using RED. For other related papers, see for instance [11] (in which a priority game is considered for competing connections sharing a drop-tail buffer), [1] as well as the survey [2]. In [15], the authors present mechanisms (e.g., AIMD of TCP) to control end-user transmission rate into differentiated services Internet through potential functions and corresponding convergence to Nash equilibrium. These references have not studied the Stackelberg equilibrium concept.

Stackelberg equilibrium has been used in other contexts of networking in [4, 9]. Both references consider M/M/1 type models for congestion. In our paper we model both TCP behavior as well as real time traffic, both sharing a common RED type router as a bottleneck. Other Stackelberg games in networking which are not directly related to TCP or to RED are [7, 18].

The structure of this paper is as follows. In Section 2 we describe the model of RED, then in Section 3 we compute the throughputs and the loss probabilities of the connections for given priorities chosen by the connections. In Section 4 we obtain the optimal priority choices of the connections for given pricing strategies of the network. The optimal pricing is then discussed in Section 5.

2 The model

RED is based on the following idea: there are two thresholds q_{\min} and q_{\max} such that the drop probability is 0 if the average queue length q is less than q_{\min} , 1 if it is above q_{\max} , and $p_{\max}(i)(q - q_{\min})/(q_{\max} - q_{\min})$ if it is q with $q_{\min} < q < q_{\max}$; the latter is the *congestion avoidance* mode of operation. This is illustrated in Figure 1.

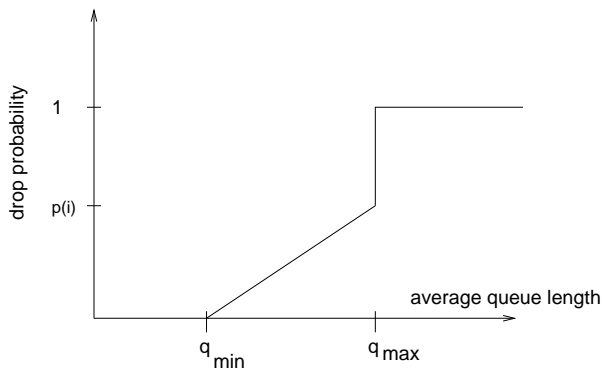


Figure 1: Drop probability in RED as function q

We consider a set \mathcal{N} containing N TCP flows (or TCP friendly flows) for data transfer and a set \mathcal{I} containing I real time flows, assumed to be TCP

friendly (in the sense that they are rate controlled according so as to achieve the same throughput that TCP achieves as a function of loss probabilities and round trip time, see eq. (1)). We assume that connections can be differentiated by means of RED algorithm; they all share a common buffer yet RED treats them differently¹. We assume that they all have common values of q_{\min} and q_{\max} but each flow i may have a different value of $p_{\max}(i)$, which is the value of the drop probability as the average queue tends to q_{\max} (from the left), see Figure 1. Denote $\mathbf{p} = \{p_{\max}(i), i \in \mathcal{I} \cup \mathcal{N}\}$. We identify $p_{\max}(i)$ as the priority class of a connection. The service rate of the bottleneck router is given by μ .

3 Computing the throughputs

We use the well-known relation for TCP rate which characterizes both data as well as real time connections (assumed to be TCP friendly):

$$\lambda_i = \frac{1}{R_i} \sqrt{\frac{\theta}{p_i}}, \quad i \in \mathcal{N} \cup \mathcal{I}, \quad (1)$$

where R_i and p_i are TCP flow i 's round trip time and drop probability, respectively. Parameter θ is typically taken as $3/2$ (when the delayed ack option is disabled) or $3/4$ (when it is enabled).

We assume that TCP senders and the rates of the TCP-friendly real time connections are not limited by the receiver window.

We model the bottleneck as a fluid queue. We assume that the buffer size is well dimensioned, i.e. it is sufficiently large so that full utilisation of the service rate can be achieved. This gives

$$\sum_{j \in \mathcal{I} \cup \mathcal{N}} \lambda_j (1 - p_j) = \mu. \quad (2)$$

The parameters λ_j and p_j are assumed to be determined by an agreement with user j . More precisely,

- p_j 's are assumed to be controlled at the RED router; each j , p_j is assumed to be a function of the priority level of connection j which is determined by user j and which we identify with the price x_j per packet payed by the j th connection.
- If μ were fixed, the p_j would then determine λ_j through (1).
- We can thus view equation (2) as providing us with the rate μ that should be available at the bottleneck link so as to guarantee the performance measures (loss probabilities and throughputs) payed for by the connections.

¹RED punishes aggressive flows more by dropping more packets from those flows

Next, we define some pricing strategy that gives p_i in terms of x_i . From Figure 1 for each i the drop probability is

$$p_i = p_{\max}(i)Q, \quad \forall i, \quad \text{where} \quad Q = \frac{q - q_{\min}}{q_{\max} - q_{\min}}. \quad (3)$$

We assume, as already mentioned in the introduction, that q is a fixed target value of the average queue. Note that $0 < Q < 1$. We assume that the cost per packet x_i is inversely proportional to $p_{\max}(i)$. Thus $p_{\max}(i)$ is given in terms of x_i by $(\alpha x_i + \beta)^{-1}$ so that

$$p_i(x_i) = \frac{Q}{\alpha x_i + \beta}. \quad (4)$$

where fixed parameters Q, α, β stand for all connections. In addition, each connection may have a fixed subscription price S , independent of the quality of service.

For the network, the main difference between data connections and real-time connections can be represented in their different utilities, i.e. in the way they perceive quality of service and prices. Data connections using TCP are not sensitive to losses since lost packets are retransmitted allowing to recover packet losses. We thus assume that their utility has a component linear in the throughput and another one related to the costs. The utility for user i ($i \in \mathcal{N}$) is thus:

$$U_i(x_i) = \lambda_i(x_i)(1 - p_i(x_i)) - a_i \lambda_i(x_i)(1 - p_i(x_i))x_i - S \quad (5)$$

where a_i is a parameter representing the weight of the price component in the utility with respect to the throughput component.

For real time connections, we assume that the utility has a component concave in the throughput, another component that represents (direct) sensitivity to losses and a component representing the evaluation of the cost for obtaining the requested priority level. We thus chose the utility function to be:

$$U_i(x_i) = w_i \log_{10} \left(1 + \lambda_i(x_i)(1 - p_i(x_i)) \right) - b_i p_i(x_i) - a_i \lambda_i(x_i)(1 - p_i(x_i))x_i - S \quad (6)$$

with w_i and b_i constant values depending on each user i .

4 Optimizing user's prices

4.1 A numerical example

We begin with a numerical example. The network parameters are given by $Q = 0.5$, $\alpha = 10$, $\beta = 200$; the parameters describing data or real time connections are $\forall i R_i = 0.05$, $\theta = 1.5$, $a_i = 0.01$. (We do not specify the value of S since it is a constant that does not affect the optimization problem faced by a user.) Furthermore the real time users are assumed to be such that $b_i = 0.1$ and $w_i = 1600 \forall i$. The functions (5) and (6) are depicted in Figures 2-3 as a function of the price x (equivalent to giving the priority level). Their restriction

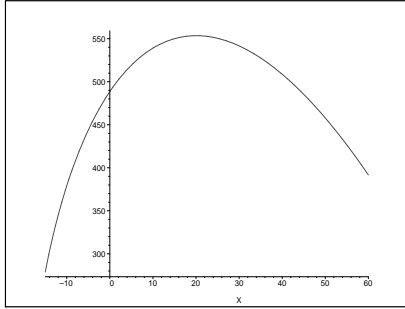


Figure 2: Utility as a function of the priority for data traffic

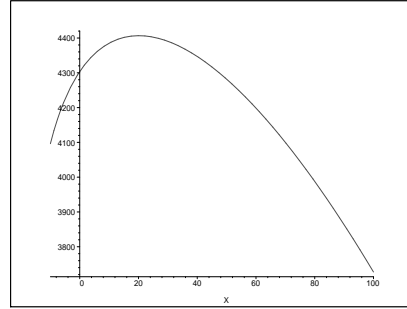


Figure 3: Utility as a function of the priority for real time traffic

to nonnegative x provides the utilities. We see that the utility as a function of the priority level x has a clear unique maximum. We shall prove this observation in the next subsection for general choice of parameters.

Below we address the question of how user i should select a priority level x so that the resulting parameters $\lambda_i(x)$ and $p_i(x)$ maximize its utility. This will determine the reaction of users to a given pricing strategy of the network. In the next section we then introduce the network's utility which will be used to compute the network's pricing strategy that maximizes its utility when taking into account the reaction of the users to that strategy.

4.2 Analysis of the user's optimization problem

Real-time connections Expression (5) leads us to the following proposition on users' best responses.

Proposition 1. *There exists a unique maximum of $U_i(x_i)$ at $X_i > 0$ for data connections. If*

$$a_i < \frac{\alpha \beta + Q}{2\beta \beta - Q},$$

it is the unique solution of given by

$$X_i = \left\{ x : \left. \frac{dU_i(x_i)}{dx_i} \right|_x = 0 \right\}$$

i.e.,

$$X_i = \frac{1}{\alpha 6} \left(\frac{\alpha}{a_i} - 5\beta + Q + \sqrt{\left(\frac{\alpha}{a_i} + \beta + Q \right)^2 + 12Q \left(\frac{\alpha}{a_i} + \beta \right)} \right). \quad (7)$$

Otherwise the maximum is located at $X_i = 0$.

Proof: The derivative of $U_i(x_i)$ is given by

$$\frac{dU_i}{dx_i}(x_i) = \left[\sqrt{\alpha x_i + \beta} + \frac{Q}{\sqrt{\alpha x_i + \beta}} \right] \times \left(\frac{\alpha(1 - a_i x_i)}{2(\alpha x_i + \beta)} - a_i \right) + \frac{2a_i Q}{\sqrt{\alpha x_i + \beta}} \quad (8)$$

$$\left[= H(X_i, a_i) \right].$$

The unique solution X_i to $\frac{dU_i}{dx_i}(x_i) = H(X_i, a_i) = 0$ is given by (7).

We have $\lim_{x_i \rightarrow +\infty} \frac{dU_i}{dx_i}(x_i) \rightarrow -\infty$; straightforward calculations show that $\frac{dU_i}{dx_i}(0) > 0$ if, and only if, $a_i < \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$, thus the solution X_i corresponds to an absolute maximum of the function $U_i(x_i)$ [17]. Otherwise, if $\frac{dU_i}{dx_i}(0) \leq 0$ we have an absolute maximum at $X_i = 0$ because the function $U_i(x_i)$ is strictly decreasing for all $x_i \geq 0$. \square

Remark that, using this proposition, the drop probability $p_i(X_i)$ is

$$p_i(X_i) = \frac{6Q}{\left(\varphi_i + Q + \sqrt{(\varphi_i + Q)^2 + 12Q\varphi_i} \right)}$$

if $a_i < \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$, where $\varphi_i = \frac{\alpha}{a_i} + \beta$, and $p(0) = Q/\beta$ otherwise.

Real-time connections Similarly the approach for the data connections, one can get the optimal choice x_i for real-time users who maximize their utilities $U_i(x_i)$.

To illustrate Proposition 1 we consider the following example. the parameter values $Q = 0.5$, $\alpha = 10$, $\beta = 20$, $\theta = 1$, $R_i = 0.05$, $a_i = 0.1 \forall i$. The value of $\frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$ results in $0.2628 > 0.1 = a_i$. Figure 4 shows the utility $U_i(x)$ of a single data connection as a function of x for these parameters. The maximum is then obtained at $X_i = 2.06 > 0$ as expected.

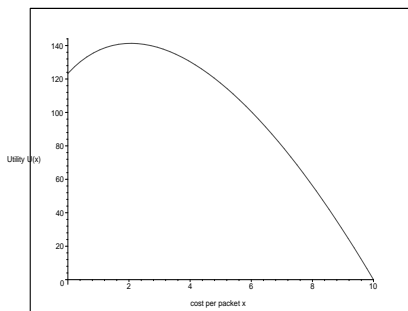


Figure 4: Utility as a function of x

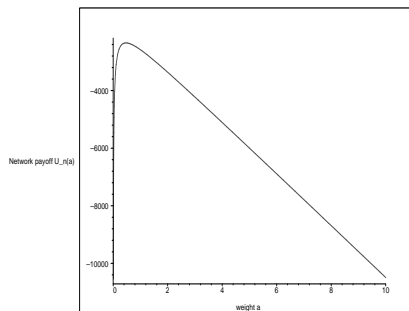


Figure 5: Network payoff as a function of a

5 Optimizing the network's pricing strategy

We assume that the network's payoff is made of two components: a cost which is proportional to the resources it provides, that is to the link rate μ , and a

utility that represents its revenues. The latter is assumed to be simply the sum of prices payed by the users. Considering a Stackelberg approach, the network's utility is thus given by:

$$\begin{aligned}\mathcal{U}_{\text{network}}(a_i) &= \sum_{i \in \mathcal{I} \cup \mathcal{N}} (a_i \lambda_i(X_i)(1 - p_i(X_i))X_i + S) - \delta \mu \\ &= \sum_{i \in \mathcal{I} \cup \mathcal{N}} (a_i \lambda_i(X_i)(1 - p_i(X_i))X_i + S) - \delta \sum_{i \in \mathcal{I} \cup \mathcal{N}} \lambda_i(X_i)(1 - p_i(X_i))\end{aligned}$$

i.e., it depends on the equilibrium obtained at users' level.

Identical TCP senders Consider N identical TCP senders, i.e., $X_i = X = X(a) > 0$, given by (7), $a_i = a$, $\lambda_i = \lambda$, and $R_i = R$. Network utility $\mathcal{U}_{\text{network}}$ simplifies to

$$\mathcal{U}_{\text{network}}(a) = N(aX(a) - \delta)\lambda(X(a))(1 - p(X(a))) + N \cdot S \quad (9)$$

As before, the maximum of (9) is obtained by using equation $d\mathcal{U}_{\text{network}}(a)/da = 0$. The solution to this expression depends also on δ . In particular, for $\delta = 1$, $\mathcal{U}_{\text{network}}(a)$ has a maximum $a = \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q} > 0$. For $\delta > 1$ then it yields

$a > \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$. In these two cases, i.e., when the optimal pricing parameter a ,

that the network imposes, satisfies $a \geq \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$, the utility of the users has

an absolute maximum at $X = 0$, see Proposition 1. This means that the users subscribe to the minimal quality of service: they pay only the subscription fee S and obtain the largest drop probability, Q/β (and hence the minimal throughput). Consider a numerical example, illustrated on Figure 5: a plot of $U(a)$ is given for $N = 10$ identical TCP connections, for $\delta = 2$, $Q = 20/40 = 0.5$, $\alpha = 10$, $\beta = 20$, $c = 1$, $R = 0.05$, $a = 0.1$. The maximum then is obtained at $a = 0.4708188006 > 0.2628$.

At $0 < D < \delta < 1$, it yields a reasonable value $0 < a < \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$; the quantity D is such that the absolute maximum is obtain at $a = 0$, on the interval $[0, \infty)$. In Figure 6a we present a plot for $\delta = 35/100$, where $a = 0.01255526875 < 0.2628$. We have $D = 0.3333$. In this case we obtain $X = 25.28249042$, the cost per packet for the maximum user's utility, see Figure 6b.

6 Conclusions

In this paper, we have analyzed a Stackelberg game where users choose the price they are willing to pay at a RED buffer in order to discriminate service. From the resulting equilibrium that we compute, the network With respect to a previous result, we use the fact that there is a fixed target for the average queue length. This simplification allows to obtain explicit formulae for the equilibrium

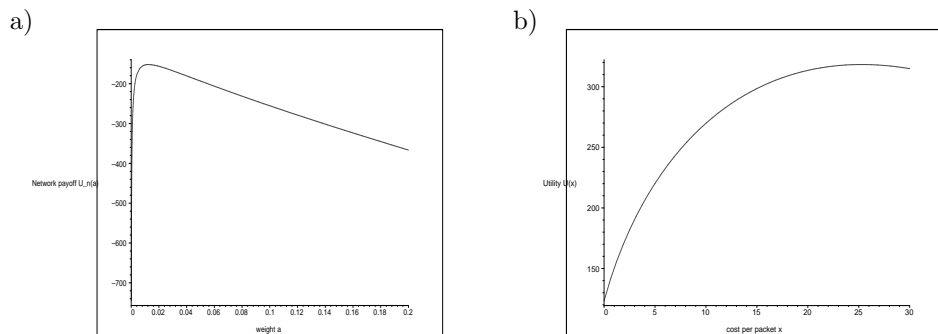


Figure 6: Network payoff as a function of a . Case $a < \frac{\alpha}{2\beta} \frac{\beta + Q}{\beta - Q}$

priority choices of users. Thanks to those expressions, we are able to solve the Stackelberg formulation of the problem.

References

- [1] T. Alpcan and T. Basar, “A game-theoretic framework for congestion control in a general topology networks”, *41st IEEE Conference on Decision and Control*, Las Vegas, Nevada, Dec. 10-13, 2002.
- [2] E. Altman, T. Boulogne, R. El Azouzi, T. Jimenez and L. Wynter, “A survey on networking games”, *Computers and Operations Research*, Vol. 33, Issue 2, 286–311, 2006.
- [3] E. Altman, R. El-Azouzi, D. Barman, D. Ross and B. Tuffin, “Pricing Differential Services: A Game-Theoretic Approach”, *Computer Networks*, 50(7), pp. 982–1002, 2006.
- [4] T. Basar and R. Srikant, “A Stackelberg network game with a large number of followers”, *J. Optimization Theory and Applications*, 115(3):479-490, December 2002
- [5] F. Bernstein and A. Federgruen, “A general equilibrium model for decentralized supply chains with price- and service-competition”, Available at <http://faculty.fuqua.duke.edu/~fernando/bio/>
- [6] D. Dutta, A. Goel and J. Heidemann, “Oblivious AQM and Nash Equilibria”, *IEEE Infocom*, 2003.
- [7] Rahul Garg, Abhinav Kamra, Varun Khurana, ”A game-theoretic approach towards congestion control in communication networks”, *ACM SIGCOMM Computer Communication Review*, Volume 32 Issue 3, July 2002.

- [8] J. Janssen, D. De Vleeschauwer, M. Büchli and G. H. Petit, “Assessing voice quality in packet-based telephony”, *IEEE Internet Computing*, pp. 48–56, May–June, 2002.
- [9] Y. A. Korilis, A. A. Lazar and A. Orda, “Achieving network optima using Stackelberg routing strategies”, *IEEE/ACM Transactions on Networking*, 5(1), pp. 161–173, 1997.
- [10] T.V. Lakshman and U. Madhow, “The performance of TCP/IP for networks with high bandwidth-delay products and random loss”, *IEEE/ACM Transactions on Networking*, Jun 1997.
- [11] M. Mandjes, “Pricing strategies under heterogeneous service requirements”, *Computer Networks* **42**, pp. 231–249, 2003.
- [12] M. Mathis, J. Semke, J. Mahdavi, and T. Ott, “The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm”, *ACM Computer Communication Review*, Jul 1997.
- [13] V. Misra, W.-B. Gong, and D. Towsley, “Stochastic differential equation modeling and analysis of TCP-window size behaviour”, *Performance*, Oct 1999.
- [14] P. Piedad, J. Ethridge, M. Baines and F. Shallwani, *A Network Simulator Differentiated Services Implementation*, Open IP, Nortel Networks, July, 2000. Available at <http://www.isi.edu/nsnam/ns>
- [15] Y. Jin, G. Kesidis, “Nash equilibria of a generic networking game with applications to circuit-switched networks”, *IEEE INFOCOM '03*
- [16] Numerical Recipes in C, The Art of Scientific Computing, 2nd Edition, Section 5.6 <http://www.ulib.org/webRoot/Books/NumericalRecipes/bookc.html1>
- [17] W. Rudin, *Functional Analysis*, McGraw-Hill, New York, 1973
- [18] Srinivas Shakkottai, R. Srikant, ”Economics of network pricing with multiple ISPs”, *IEEE/ACM Transactions on Networking (TON)*, Volume 14 Issue 6, December 2006.