

Network control and usage-based charging : Is charging for volume adequate?*

C. Courcoubetis[†], V. A. Siris, and G. D. Stamoulis[†]

Institute of Computer Science (ICS)

Foundation for Research and Technology - Hellas (FORTH)

P.O. Box 1385, GR 711 10 Heraklion, Crete, Greece

Tel: +30 81 39 1726, Fax: +30 81 39 1601

{courcou, vsiris, gstamoul}@ics.forth.gr

Abstract

There has been extensive research on the application of the effective bandwidth concept for quantifying resource usage in order to create simple yet effective usage-based charging schemes with desirable incentive properties. This research has shown that simple charging schemes which involve two measurements, time and volume, can serve their purpose well by producing adequate approximations of the effective usage of a bursty traffic stream. An issue that has not been addressed in detail is the relation between the type of admission control mechanism the network uses and the definition of effective usage. In particular, the above charging schemes assume "static" Connection Admission Control (CAC). In contrast to static CAC, "dynamic" CAC strategies utilize on-line measurements of the actual load, hence can achieve much higher utilization. We argue that under such dynamic strategies the effective usage concept must be redefined, and that when control actions occur at faster time scales than the burstiness of the sources, the effective usage approaches the mean rate of the sources. In addition to justifying the potential use of simple volume-based tariffs, the above justifies the deployment of sophisticated dynamic admission control mechanisms, since these result in more competitive prices.

Keywords: usage-based charging, effective bandwidths, incentives, time scales, Connection Admission Control

1 Introduction and overview

The increasingly competitive nature of the telecommunications market and steps in deregulation are pushing towards prices which take some account of actual resource usage. Usage-based pricing enables network providers to recover the costs of service provision from their customers in a fair manner. Of course the price for network services will be

*This work was supported in part by the European Commission under ACTS Project CASHMAN (AC-039).

[†]Also with the Dept. of Computer Science, University of Crete.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICE 98 Charleston SC USA

Copyright ACM 1998 1-58113-076-7/98/10...\$5.00

affected, to a large degree, by competition and marketing issues. However, it is unlikely that a pricing scheme that does not take into account resource usage is an effective mechanism for controlling congestion and for providing the right incentives for efficient and stable network operation. The design of simple tariffs that do not sacrifice incentive compatibility, i.e., tariffs that provide the incentives for cost-minimizing users to select traffic contracts and use the network in ways that maximizes the aggregate utility of all the users, is clearly an important issue. An example of simple tariffs is charging only for volume; it has been argued that such tariffs, due to their simplicity, lack to a large extent incentive compatibility properties. In this paper we argue that charging for volume can be adequate under reasonable assumptions on the underlying network control technology and the time scales of the burstiness of the traffic sources.

We start our discussion by motivating the use of tariffs that account for resource usage. A prime example to consider is that of the Internet, which faces intense congestion problems. Many engineers and economists believe that the Internet's congestion problems are due, at least in part, to its ineffective pricing structure, namely *flat rate* pricing, where prices depend only on the rate of the access pipe which connects a customer to his provider [14, 6]. Flat rate pricing reveals mostly user preferences for connectivity rather than quality, and hence results in highly congested networks offering connectivity services (email, ftp, etc.) at very low costs. Such a pricing scheme provides no incentives for users to use less bandwidth than the rate of their access pipe. Furthermore, flat rate pricing does not enable users to adequately reveal their preferences for network usage. All users are treated the same, even though different users might have a different value for the same service. Both of these limitations result in a congested network where resources are not used according to the actual user needs.

Unlike traditional telephone networks, which offer a single service that occupies a constant amount of bandwidth, broadband networks will support connections with different traffic characteristics and Quality of Service (QoS) requirements, expressed in terms of loss probability and delay. In addition to services offering some level of QoS guarantees, broadband networks will offer services targeted for elastic traffic [15]. For such services, the information transfer capability offered to users is no longer statically defined at call

setup, but can change throughout the lifetime of a connection. Pricing and resource sharing for such services is related to flow control, and an important issue is the appropriate definition of fairness, given the constraints of the time scales of the feedback loop of flow control. The interested reader is referred to [12, 13, 4].

In this paper we focus solely on services where a traffic contract is statically defined at call setup. This contract specifies the maximum traffic a user is allowed to send into the network, and the performance that the network will guarantee. Note that this category includes not only the Constant Bit Rate (CBR) and Variable Bit Rate (VBR) services in ATM networks, and the guaranteed and controlled-load services in the Internet's integrated services architecture, but also the connection of a large user with his Internet service provider.

Clearly the quality of the network service and the statistical properties of traffic should be reflected on a user's charge since the network, in order to guarantee its QoS, must reserve a corresponding amount of resources. Hence, an important first step to charging is to find a workable and accurate way for quantifying resource usage which takes into account the above characteristics. The *effective bandwidth*, which is based on rigorous mathematical justification [11], is a scalar which summarizes the amount of resources required by a connection in order to preserve its own QoS requirements, and the requirements of the other connections it is multiplexed with.

Recently [2], a mathematical theory for creating usage-based charging schemes using the notion of the effective bandwidth has been developed. The approach, which builds on the initial work of [10], allows for the design of tariffs that can take into account an arbitrary number of measured traffic parameters. Experiments with real traffic (Internet WAN and MPEG-1 compressed video) [1] have shown that an important class of the above tariffs, where charges are based solely on time and volume, presents a good trade-off between accuracy and accounting overhead, hence can serve their purpose well. Even though being simple, these charging schemes preserve to a large extent their incentive compatibility properties and are natural to implement with current technology.

An issue that has not been addressed in detail is related to the time window over which resource usage is measured. For example, resource usage need not be measured for the whole duration of a connection, but can be measured over a smaller time window. Indeed, as this time window becomes smaller than the typical time it takes for the rate of the source to change, the source's rate during such an interval tends to remain rather constant, and its relative resource usage becomes proportional to the above value. This in turn implies that the effective usage (rate) of a connection over a large number of such intervals approaches the mean rate of the source. We will argue in this paper that if the network's admission control mechanism uses on-line measurements of the actual load and results in effectively loading the system, then it is sensible to define the above time window to be the time between consecutive control decisions (accept or reject an incoming connection). This suggests a new definition for the effective bandwidth that takes into account the

time scales of the admission control mechanism. Using the previous remarks, when the time scale of the arrivals of new connections is smaller than the time scale of the burstiness of the input traffic, the effective bandwidth of a connection is well approximated by its mean rate and the relative resource usage of the connection is proportional to the total volume carried.

The rest of the paper is organized as follows. In Section 2 we review the basic results from the theory of effective bandwidths and discuss how simple usage-based charging schemes with desirable incentive properties that require two measurements, time and volume, can be engineered. Section 3 deals with Connection Admission Control (CAC), and discusses the relation of static and dynamic (or measurement-based) CAC with resource usage. We argue that in the case of dynamic CAC, the right effective bandwidth concept is different than the one originally used for static CAC. Next, Section 4 discusses how three important time scales, namely the time scale of buffer overflow, the time scale of network control, and the time scale of traffic burstiness, affect resource usage and the relative performance of dynamic and static CAC. We also discuss how these affect the form of the charge. Finally, Section 5 provides some conclusions.

2 Charging schemes linear in time and volume based on effective bandwidths

In this section we briefly review the theory of effective bandwidths, and discuss how simple time- and volume-based charging schemes, with desirable incentive properties, can be created based on the effective bandwidth.¹

2.1 Effective bandwidths of stationary sources

If $X_j[0, t]$ is the amount of workload produced by a stationary source of type j in a time interval of length t , the effective bandwidth for that type of source is defined as [11]

$$\alpha_j(s, t) = \frac{1}{st} \log E \left[e^{sX_j[0, t]} \right], \quad (1)$$

where s, t are system parameters which depend on the characteristics of the aggregate multiplexed traffic, the QoS (cell loss), and the link parameters (capacity and buffer size) [2]. Specifically, the physical interpretation of the *time* parameter t (measured in, e.g., milliseconds) corresponds to the most probable duration of the busy periods of the buffer prior to overflow. The *space* parameter s (measured in, e.g., kb^{-1}) is affected by the degree of multiplexing and depends, among others, on the size of the peak rate of the multiplexed sources relative to the link capacity. In particular, for links with capacity much larger than the peak rate of the multiplexed sources, s tends to zero and $\alpha_j(s, t)$ approaches the mean rate of the source, while for links with capacity not much larger than the peak rate of the sources, s is large and $\alpha_j(s, t)$ approaches the peak rate.

According to the theory of multiplexing developed in [11, 2], given a target overflow probability $e^{-\gamma}$, the constraint on the sum of the effective bandwidths of stationary

¹For a more rigorous mathematical treatment of the subject, the reader is referred to [11, 2].

sources entering a link with capacity C and buffer B has the following linear form:

$$\sum_j N_j \alpha_j(s, t) \leq C + \frac{1}{t} \left(B - \frac{\gamma}{s} \right) = C^*, \quad (2)$$

where C^* is the amount of "effective resource (capacity)", N_j is the number of sources of type j , and (s, t) is an extremizing pair of

$$\sup_t \inf_s \left[st \sum_i N_i \alpha_i(s, t) - s(Ct + B) \right]. \quad (3)$$

The optimum value of the above expression approximates the logarithm of the cell loss probability at the multiplexer.

An important property of the above definition of effective bandwidths, expressed by equation (2), is that the effective bandwidth reflects relative usage: if source A has twice as much effective bandwidth as source B (measured at a particular operating point s, t), then source A uses twice as much resources as source B, hence one source of type A can be substituted by two sources of type B. Studies in [5] have validated the above definition of the effective bandwidth and have shown that the pair (s, t) is, to a large extent, insensitive to small variations of the traffic mix and hence can be pre-computed off-line with reasonable accuracy for various periods of the day from existing traffic traces.

We conclude our discussion on effective bandwidths by stressing the importance of the time parameter t in determining the value of the effective bandwidth and hence the resource usage of a source. The definition (1) suggests that if the source changes many times during the window t (i.e., the time scales of burstiness are smaller than t), then the effective bandwidth will be close to the mean rate, while if the source changes slowly compared to t , then the effective bandwidth can be substantially larger than the mean rate.

2.2 Charging schemes resulting from effective bandwidths

One approach to charging would be to directly apply the effective bandwidth formula (1). However, this has a number of disadvantages. First, it is costly to implement since it requires for each connection traces with granularity at most t , and the computation of the logarithmic moment generating function. Second, it does not take into account the resources that the network must reserve at connection setup and does not penalize users asking for "large" traffic contracts which they will not eventually use (a user that requests at connection setup the whole link capacity but sends no traffic during the connection will pay zero). Finally, such an approach leads to a complicated charging scheme, making it difficult for users to determine the potential effect of traffic shaping on their charge.

A second approach would be to charge a connection based on an empirical estimate of the effective bandwidth for all past connections of the same type (asking for the same traffic contract). Such a charging scheme has the same disadvantage as all-you-can-eat restaurants, namely it encourages users to over-eat, and hence provides wrong incentives to users. Note that this approach resembles flat-rate pricing.

Finally, a third approach can be to charge according to the worst case effective bandwidth that could result from the given traffic contract. Such a scheme is unfair to users that happen to send less traffic than the maximum allowed by their contract, hence encourages them to send the maximum traffic that their contract allows (similar to the previous approach).

The charging approach developed initially in [10] and then in [2] takes into account both static traffic contract parameters (known a priori) and dynamic parameters (measured a posteriori). In particular, the approach shows how to transform simple tariffs of the form $aT + bV$, where T is the duration and V is the volume of a connection (measurements), into sound approximations of the effective bandwidth of the connection, by casting all the information from the static traffic contract parameters and the operating point of the network into the value of the coefficients a, b .

According to the approach, based on his traffic contract, a user is offered at connection setup a set of tariffs, a tariff being a pair (a, b) , to choose from. A rational user will select the pair (a, b) which minimizes the a priori expected value of his charge. Based on the theory developed in [10, 2], the tariff pairs (a, b) can be appropriately defined so that the expected charge for a rational user is $\bar{\alpha}T$, where $\bar{\alpha}$ is an approximation of the largest effective bandwidth that is consistent with both the traffic contract and the actual traffic measurements. In [1] we investigated, for real traffic (Internet WAN and MPEG-1 compressed video), various approximations $\bar{\alpha}$ for the effective bandwidth. The experiments have shown that the above simple time- and volume-based charging scheme results in charges that accurately reflect relative resource usage.

3 Connection admission control and effective bandwidths

In this section we define the basic concepts of Connection Admission Control (CAC) and discuss their relation to the definition of effective bandwidths. We will argue that the original definition of an effective bandwidth is sound only in the case of static control, whereas a new definition must be used for systems using dynamic control, which takes account of the available information and the granularity of the control actions.

Connections that request network services with some quality of service guarantees are subject to admission control. A simple definition of such a mechanism is the following. A connection request to the network contains the description of the user-network contract that must be honored if the network decides to accept the connection. This contract contains the minimum quality of service the network must provide and a maximum value for certain parameters (traffic contract parameters) of the traffic that will be generated by the source (possibly also the tariff that the user prefers). The network uses this information to check if it has enough resources to satisfy the new connection simultaneously with the other connections that have already been accepted. If it does, then it accepts the connection, otherwise the connection is rejected.

If one thinks of the network as a producer of services (provision of connections for carrying user traffic), then the

connection admission control mechanism defines the technology set of the network. The operating point of the network will correspond to the solution of a complex optimization problem by the network operator, which encodes other important information (demand elasticity, etc.). When the technology set is locally characterized by linear constraints of the form (2), it follows that at a given operating point the price of a connection will be proportional to its effective bandwidth and the shadow price of the effective bandwidth constraint. Hence, it is important to understand how different types of admission control mechanisms might influence the above constraints, which in turn will be reflected on the structure of the prices.

3.1 Static versus dynamic control and effective bandwidths

There are two approaches to Connection Admission Control (CAC): static and dynamic. With static CAC², the admission decision is based solely on the (static) traffic contract parameters and possibly on other a priori available information about the statistical properties and the burstiness of both the new connection and the other connections that have already been accepted. The intuition behind static CAC is that once a connection is accepted, the network, based on the traffic contract, reserves some amount of resources for the whole duration of the connection. This amount must correspond (see (2)) to the effective bandwidth of the connection as defined in the previous section (for simplicity assume that complete statistical characterization of the new source is known³), since *the admission decision makes no assumption on future connection departures (connection duration) and has no other information besides the statistical description of the sources.*

Dynamic CAC can be an extremely intelligent and complex mechanism to describe, and there are many such mechanisms proposed in the literature (e.g., see [7, 3, 9, 8] and the references therein). These, among other, differ in the measurements required and in the logic of accepting/rejecting connections. There are two crucial differences with static CAC:

- The decision on accepting/rejecting a connection is based on the *actual* traffic load and on the rate of termination of the active connections (rate of load decrease) and arrival of new connections.
- Since the decision algorithm takes risks based on assumptions of the future of the connection arrival and departure process, the QoS is defined as the average QoS over the potential states of the system (less strict than in the case of static CAC).

An important concept that is useful to define for dynamic CAC is the *conditional effective bandwidth* of a source over a control interval T . Intuitively, this corresponds to the case

²Static CAC can be used for providing both deterministic and statistical QoS guarantees. On the other hand, dynamic or measurement-based CAC can provide only statistical QoS guarantees.

³If only traffic contract information is used for CAC, because traffic descriptors, such as the leaky bucket, provide only a crude characterization of a traffic stream, static CAC can result in very low utilization.

where we multiplex a number of sources for which we have some information about their state at time 0, and we want the cell loss probability over the following T time units to be at most $e^{-\gamma}$. Simple mathematical arguments suggest that in this case a similar equation to (2) holds, with the effective bandwidth defined by (1) where the expectation is the conditional expectation given the information at time 0, and the random process is the truncation of the original one over a window of length T ⁴. For an on-off source, simple algebra shows that if the time to overflow the buffer t is less than T and T is small compared to the time it takes for the rate of the source to change, then the conditional effective bandwidth is close to the initial value of the rate (a bit less (more) if the source was on (off) initially).

A generic model for dynamic CAC is now the following. Let the control interval T correspond to the interarrival time of connection requests. Then, at the beginning of each such interval, the controller computes the sum of the conditional effective bandwidths of the accepted sources; this is the current effective load until the next decision epoch. Based on the spare capacity available, it then decides to accept or reject the new connection. Clearly a *necessary* condition is that the total effective load during the current control interval is less than the available capacity. But this is not sufficient for ensuring such a condition in the future, since the above measure of the effective load could be increasing due to the behavior of the sources, even though no new connection is accepted. Basically, the controller must assess the risk of accepting a particular effective load given the current information about the state of the sources; this is the probability where, without accepting new connections and only letting connections depart, eventually the effective load will exceed the capacity. Hence, a successful strategy will operate in a controlled risk fashion, where the carried load is maximized while keeping the risk below a certain level.

We can now argue that, if many sources are multiplexed and connection arrivals occur frequently enough, any efficient dynamic control strategy will succeed in keeping the effective load close to capacity, for an overwhelmingly large proportion of control intervals. Hence, the relative resource usage of a source is the average of the values of the conditional effective bandwidth of the source over all control intervals that occurred during its lifetime. Following the example of the on-off sources, when the time to fill the buffer is small compared to the control interval (connection interarrival time), and the latter time is small compared to the time scale of burstiness of the source, the above average is well approximated by the mean rate.

4 Implications to charging

We have already mentioned that economic efficiency requires the design of charging schemes that are incentive compatible, hence charges should reflect actual resource usage. On the other hand, competitiveness in network service provision requires the above charges to be as low as possible. Since our charging methodology charges for effective usage, our

⁴To be rigorous, in order to make the truncated process stationary, we must construct a periodic process with period T , which coincides with the original process over the given interval of time.

charging schemes are incentive compatible. Next we discuss the competitiveness issue using the results of the previous section.

Comparing the competitiveness of static versus dynamic CAC amounts to comparing their corresponding technology sets. Clearly, the more connections a network can handle, the smaller the absolute charge can be. We have showed that the call handling capability under both types of control is described by (2), under a different interpretation of the effective bandwidths. This amounts to comparing for a given connection its effective bandwidth under each admission control mechanism. We will denote the initial effective bandwidth definition as the *unconditional effective bandwidth*.

A first remark concerns the available information. Clearly, both static and dynamic CAC define smaller effective bandwidths when the available information increases. Hence, to make the comparison more interesting we assume that both mechanisms have access to complete information regarding statistical properties of the sources and that the dynamic information is as detailed as possible.

We will show that this comparison entails different results depending on the relation of the basic time scales t (typical time for filling the buffer in the multiplexer), T_b (time scale of the changes of the rate of the source, i.e., burstiness) and T (control time scale, i.e., connection interarrival time). A first reasonable assumption is that CAC will be used for real-time traffic connections, in which case we can safely assume $t < T$. This is the case since in order to obtain small delay, which is required for real-time services, the buffer size must be small, hence the time to fill it will also be small. It has been observed [5] that for small buffer sizes, and in particular for buffer sizes giving a delay of less than 4 milliseconds (which is equivalent to a buffer of less than 1405 cells in a 155 Mbps link), the time parameter t is less than 100-200 milliseconds, which is smaller than typical connection interarrival times which are of the order of seconds.

First, consider the case where the time scales of traffic burstiness are smaller than the time scale of buffer overflow t . In this case, the effective bandwidth under both control strategies is close to the mean rate of the sources (by the law of large numbers, in a window t the average contribution of such a source is close to its mean rate), and hence both admission strategies are equally competitive. Intuitively, this is because the source mixes so fast that it is very easy to multiplex in both situations, and any on-line measurement becomes almost immediately obsolete.

If the time scale of traffic burstiness T_b is larger than the time scale of buffer overflow, then both the conditional effective bandwidth and the unconditional effective bandwidth initially increase (from being equal to the mean rate) as T_b increases, while always the conditional being the smaller of the two. Then as T_b becomes larger than the interarrival time T , the unconditional effective bandwidth continues to increase and converges to some maximum value (the on-off fluid approximation), while the conditional effective bandwidth decreases and converges to the mean rate.

The implications of the previous arguments are now straightforward. First, comparing the efficiency of dynamic ver-

sus the static CAC depends on the above time scales of the multiplexing system. The relative efficiency is inversely proportional to the ratio of the conditional and the unconditional effective bandwidths, defined for the time scales of the particular multiplexing situation. For sources varying slower than the connection arrival process, this ratio is maximized and the gain for dynamic CAC can be substantial. In many cases, dynamic CAC could accept over 50% more sources (the ratio of the effective bandwidth for typical on-off sources over their mean rate).

The second implication is that for the case of slowly varying sources, charging for volume, although being extremely simple, is incentive compatible if the network uses dynamic control strategies. Since by the previous argument dynamic CAC is the more competitive control scheme, we can assume that such mechanisms will be widely deployed in the future, and hence volume-based charging will be the tariff form used. This is further justified by assuming that future networks will have large capacities and hence will support large volumes of connection requests. Hence, the connection interarrival time will be smaller than the time scale of traffic burstiness T_b , which does not depend on the size of the network.

5 Conclusion

There has been extensive research on the use of the effective bandwidth for quantifying resource usage in order to create simple yet incentive compatible usage-based charging. The classic approach uses an effective bandwidth definition which captures effective usage when the network control strategy for admitting connections uses no dynamic measurement information about the actual load, and is solely based on a priori information about the connection traffic statistics. This definition has resulted in charging schemes which in their simplest form must account for both the duration and the volume of the connection. In this paper we argue that the use of more sophisticated connection admission procedures that are based on on-line dynamic traffic information can result in more competitive prices, and that under reasonable assumptions on system parameters, the resulting charging schemes need only measurements of volume. This is of great importance for the design of future accounting systems and for defining the structure of call detail records.

Since the above schemes are accurate when sources tend to switch behavior in time scales larger than the connection interarrival time, an interesting open question remains on providing the right incentives for users to shape their traffic in order to "slow-down" the source behavior. Interestingly enough, slow sources are easier to multiplex under dynamic connection admission strategies.

References

- [1] C. Courcoubetis, F. P. Kelly, V. A. Siris, and R. Weber. A study of simple usage-based charging schemes for broadband networks. In *Proc. of IFIP International Conference on Broadband Communications (BC'98)*, Stuttgart, Germany, April 1998.

- [2] C. Courcoubetis, F. P. Kelly, and R. Weber. Measurement-based charging in communications networks. Technical Report 1997-19, Statistical Laboratory, University of Cambridge, 1997.
- [3] C. Courcoubetis, G. Kesidis, A. Ridder, J. Walrand, and R.R. Weber. Admission control and routing in ATM networks using inferences from measured buffer occupancy. *IEEE Trans. Commun.*, 43:1778-1784, 1995.
- [4] C. Courcoubetis and V. A. Siris. An approach to pricing and resource sharing for Available Bit Rate (ABR) services. Technical Report No. 212, ICS-FORTH, November 1997.
- [5] C. Courcoubetis, V. A. Siris, and G. D. Stamoulis. Application and evaluation of large deviation techniques for traffic engineering in broadband networks. In *Proc. of ACM SIGMETRICS'98/PERFORMANCE'98*, Madison, Wisconsin, June 1998.
- [6] E. Firdman. Rx for the Internet: Usage-based pricing. *Data Communications*, January 1997. Available at http://www.data.com/business_case/rx_internet.html.
- [7] R. J. Gibbens, F. P. Kelly, and P. B. Key. A decision-theoretic approach to call admission control in atm networks. *IEEE J. Select. Areas Commun.*, 13(6):1101-1114, August 1995.
- [8] M. Grossglausser and D. Tse. A framework for robust measurement-based admission control. In *Proc. of ACM SIGCOMM'97*, Cannes, France, September 1997.
- [9] S. Jamin, S. J. Shenker, and P. B. Danzig. Comparison of measurement-based admission control algorithms for controlled-load service. In *Proc. of IEEE INFOCOM'97*, Kobe, Japan, April 1997.
- [10] F. P. Kelly. On tariffs, policing and admission control for multiservice networks. *Operations Research Letters*, 15:1-9, 1994.
- [11] F. P. Kelly. Notes on effective bandwidths. In F. P. Kelly, S. Zachary, and I. Zeidins, editors, *Stochastic Networks: Theory and Applications*, pages 141-168. Oxford University Press, 1996.
- [12] F. P. Kelly. Charging and rate control for elastic traffic. *European Transactions on Telecommunications*, 8:33-37, January 1997.
- [13] F. P. Kelly, A. Maulloo, and D. Tan. Rate control in communication networks: shadow prices, proportional fairness and stability. *Journal of the Operational Research Society*, 49, 1998.
- [14] J. K. Mackie-Mason and H. R. Varian. Pricing the Internet. In B. Kahin and J. Keller, editors, *Public Access to the Internet*. Prentice Hall, Englewood Cliffs, New Jersey, 1995.
- [15] S. Shenker. Service models and pricing policies for an integrated services Internet. In B. Kahin and J. Keller, editors, *Public Access to the Internet*. Prentice Hall, Englewood Cliffs, New Jersey, 1995.