



**HAL**  
open science

# Charging the Internet without Bandwidth Reservation: an Overview and Bibliography of Mathematical Approaches

Bruno Tuffin

► **To cite this version:**

Bruno Tuffin. Charging the Internet without Bandwidth Reservation: an Overview and Bibliography of Mathematical Approaches. [Research Report] RR-4355, INRIA. 2002. inria-00072233

**HAL Id: inria-00072233**

**<https://inria.hal.science/inria-00072233v1>**

Submitted on 23 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

*Charging the Internet without bandwidth  
reservation: an overview and bibliography of  
mathematical approaches*

Bruno Tuffin

**N°4355**

Janvier 2002

———— THÈME 1 ————



*Rapport  
de recherche*



# Charging the Internet without bandwidth reservation: an overview and bibliography of mathematical approaches

Bruno Tuffin \*

Thème 1 — Réseaux et systèmes  
Projet Armor

Rapport de recherche n° 4355 — Janvier 2002 — 32 pages

**Abstract:** Pricing is one of the biggest challenges of the next generation of the Internet. Even if flat rate pricing is probably one of the main reasons of the Internet success, the only way to prevent from network congestion and to differentiate services is by means of usage-based pricing schemes. We review in this paper, from a mathematical modeling point of view, the pricing schemes *without resource reservation* that have been developed in the literature. Indeed, an advantage of the absence of reservation in the Internet is that network management is cheap. Even if accounting and billing will increase this cost, we believe that pricing without resource reservation is the lesser of two evils with respect to applying some costly bandwidth reservation procedures.

**Key-words:** Fairness, Internet economics, Optimization, Pricing, Service Differentiation.

(Résumé : *tsvp*)

\* [btuffin@irisa.fr](mailto:btuffin@irisa.fr)

# Tarification de l'Internet sans réservation de bande passante : un état de l'art et une bibliographie des approches mathématiques

**Résumé :** La tarification est l'un des plus grands défis de la prochaine génération de l'Internet. Même si la tarification basée sur une utilisation illimitée est probablement une des raisons principales du succès de l'Internet, le seul moyen de prévenir l'engorgement du réseau et de différencier les services est d'utiliser des modes de tarification basés sur l'utilisation. Nous listons dans cet article, d'un point de vue mathématique, les schémas de tarification *sans réservation de ressources* qui ont été développés dans la littérature. En effet, un avantage de l'absence de réservation dans l'Internet est que la gestion du réseau est peu coûteuse. Même si mesurer et facturer augmenteront ce coût, nous estimons que la tarification sans réservation de ressources est un moindre mal par rapport à une très coûteuse réservation de bande passante.

**Mots-clé :** Différenciation de services, Économie de l'Internet, Équité, Optimisation,

## 1 Introduction

The Internet is undergoing a tremendous development of its traffic. A consequence is that real users complain that large data transfers take too long, without any possibility to improve this by themselves (by paying more for instance). To cope with this congestion, it is possible to develop the link capacities but many authors consider that it is not a viable solution as the network must respond to increasing demand (and experience has shown that demand of bandwidth has always been ahead of supply), especially now that the Internet is becoming a commercial network. Furthermore, incentives to a fair utilization between customers is not included in the current Internet (see for instance [33, 45]). For these reasons, it is suggested that the current flat rate fees, where customers pay a subscription and obtain an unlimited usage, are replaced by usage-based fees [13]. Also, the future Internet will supply different kinds of services such as video, voice, email, ftp, telnet or html among others. Each of these applications requires a different quality of service (QoS): for example, video needs very small delays and packet losses, voice requires small delays but can afford some cell losses, email can afford delay (within a given bound) while ftp needs more a good average throughput and telnet is more interested in small round trip times. Some pricing incentives should exist so that each user does not always choose the best QoS for his application and so that the final result is a fair utilisation of the bandwidth. On another hand, we need to be aware of the trade-off between engineering efficiency and economic efficiency; indeed, measurement for example allows to improve the management of the network but is costly.

In [66], J. Roberts classifies pricing schemes in three categories, flat rate pricing, congestion pricing and transaction pricing and studies their impact on QoS (see also other introductory or overview papers [12, 13, 14, 29, 74]; and [64] where a interesting time-scale methodology and classification is presented). Another classification separates the schemes between edge pricing, where the charge is set only at the edge of the network, and node-per-node pricing. Our paper differs from the previous ones in that mathematical models are displayed (when available). We classify the suggestions about the future Internet pricing in eight different families as follows, most of them being sub-categories of congestion pricing in [66].

1. As already explained, a first group of people (see for instance [1, 59]) are arguing that even if the number of customers (and their demands) is growing quickly, the network capacity is adapting itself to the demand. Furthermore, if the system has survived so far and has known such a success, why should we introduce a costly billing model?
2. For a second group of people, an incentive pricing will be necessary to regulate these various quality of services and some services must be *guaranteed*. Like in ATM networks, charging models for guaranteed services such as voice or video should be related to connection acceptance control (CAC) [15, 18, 19, 35, 73], resource reservation and effective bandwidth theory [31]. In the Internet the reservation of resources may be done using RSVP [78]. In [76] and [21], reservation is used and only non guaranteed services are accomplished using best effort techniques adapted to the willingness to pay of the user. In [61], CAC and bandwidth reservation is applied to loss networks; a nice characteristic is that arrival rates for each class of service depend on the connection fee of the class. A dynamic programming method is used to obtain optimal and quasi-optimal prices and it is shown that time-of-day pricing is efficiently approximating congestion pricing. These results are extended in [44] to general loss networks (non-exponential holding times) and to the case where the system has a prior knowledge of connection times, at their arrival. In [56], pricing elastic traffic flows is related routing.
3. Another alternative has been suggested by A. Odlyzko in [58]. The proposal is called Paris Metro Pricing (PMP) by analogy with the Paris Metro System. The network is decomposed into several separate networks and each network, working like the current Internet, has a different connection fee so that we expect that the most expensive ones will not be less congested. Thus no QoS is guaranteed but the model can be easily implemented without huge overhead.
4. The Cumulus Pricing Scheme (CPS) [62, 63] is also a simple possibility. A contract is negotiated between the ISP and the user. During periods of time, the utilisation is measured and (positive or negative) cumulus

points are awarded, depending on whether the contract is satisfied or not. At a given time extra-fees can be charged.

5. Another group suggests to use priority pricing, without reservation of resources (see [5, 8, 9, 25, 30, 55] and the references therein). Each class is assigned a priority number and is served according to this policy at each node of the network. Priority pricing scheme is decomposed into two sub-classes:
  - (a) the first one is posted priority pricing where each priority class price is established in advance. In [5], each customer is assigned a quota for high priority packets (following his contract) and if his quota is exceeded, he has a penalty in charge the next month. In [8], a priority flag is assigned to each packet according to the type of service, but also a reject flag for services which can bear some losses. In [49, 50, 48], a discrete time model is described where the time is divided in time slots. Optimal prices are computed in order to maximize the network benefits.
  - (b) The second sub-class is non-posted priority pricing where the price of the packet class depends on the traffic level. In [25], an adaptative priority pricing depending on the context (similar to the principles in [32]) is used. In [27, 55], an optimal incentive-compatible pricing scheme for the M/M/1 multi-class queue is studied (note that the result can be easily extended to the M/G/1 queue).
6. Bidding for priority has also been proposed in [53, 54]. The user makes a bid for each packet and only bids greater than some cutoff values are admitted. In [43, 68, 69], auctions for packets are replaced by auctions for bandwidth during intervals of time to reduce the management overhead. Efficiency, stability and fairness issues are solved in the case of one node but also in the case of interconnected networks.
7. Another scheme is the *expected capacity* theory developed by Clark [7] where packets are flagged *in* or *out* and are served without priority except in the case of congestion where *out* packets receive a congestion pushback.



8. A last group of pricing schemes is charging for elastic traffic based on transfer rates (see also [77]). In Kelly et al.'s work [33, 36], the users decide their payment and receive as transmission rate what the network allocates to them. In Low et al.'s work [3, 2, 46, 45, 47], the users decide their rates and pay for it according to the price computed by the network. A variant of Kelly et al.'s work has been given by La and Anantharam in [39, 40] where the flow rates are actually controlled by the window-based algorithm of TCP connections.

A vast literature has been developed in recent years on the future Internet and the integration of different services [5, 7, 8, 9, 20, 25, 28, 34, 58, 71] as well as fairness issue (see for instance [6, 51, 52, 57, 67]). To say which of the different charging groups will be implemented and how prominent they will be in practice is right now a bet. Like in [25], we believe that the arguments in favor of simply overprovisioning the capacity of the Internet is dangerous in the current status. Moreover, the capacity reservation for some types of services is expensive to implement. We are then betting that the next pricing scheme will be pricing without bandwidth reservation.

The aim of this paper is to review current works (when possible, the mathematical models) on pricing without bandwidth reservation theory (from 3 to 8 in the previous classification).

## 2 Paris Metro Pricing [58]

The proposal in [58] is to partition a network into several logically separate networks (or classes), each having a fixed fraction of the capacity of the entire network. All networks would route packets according to the current TCP and UDP. There is no formal guarantee of QoS, but by charging different rates for different classes (served in the same way), it is supposed that the most expensive classes will be less congested by self-regulation and then will deliver a better QoS. The name given to this model, *Paris Metro Pricing* (PMP) stems from the rules of the Paris Metro about 20 years ago where two class cars were existing in trains, with exactly the same quality of seats. As tickets prices were different, the cars for the most expensive class were less congested leading to a better perception of QoS.

The advantage of PMP pointed out by Odlyzko is that, even if we do not guarantee any QoS using this scheme, it would permit dispensing with measures such as RSVP and their complexity, and keep the simpler and cheaper current model of the Internet.

It is suggested that only few (3 or 4) subnetworks are implemented to minimize losses from not aggregating all the traffic. PMP charges would be assessed on each packet, and would probably consist of a fixed charge per packet and a fee depending on the size of the packet.

Recently, in [17], Gibbens et al. have studied PMP in the case of two Internet service providers (ISP) competing to maximize their profits. In their paper, a user joins the network  $i$  which maximize his utility  $U(\theta, i) = V - \theta Q^i / C^i - p^i$  where  $V$  is the positive valuation of the user,  $\theta$  is his preference for lack of congestion ( $\theta$  is assumed to follow a uniform distribution in  $[0, 1]$ ),  $Q^i / C^i$  is the mass of users divided by the capacity (at network  $i$ ), i.e., the measure of congestion, and  $p^i$  is the price per unit time charged by network  $i$ . Network  $i$  then tries to maximize its benefit  $p^i Q^i$ . It is shown that, at the stable point, the ISP will not provide multiple services. By then, they state that PMP may not survive under competition (at least if the system follows the given assumptions).

### 3 The Cumulus Pricing Scheme (CPS) [62, 63]

Like PMP, CPS is interesting by its implementation simplicity. In this scheme, the user negotiate with the ISP a given utilisation or a given QoS during a period of time. Say for instance that the contract is the volume of packets sent. If this volume is  $V(t)$  at time  $t$  and is measured between period  $[t_{i-1}, t_i]$ , the over or under-utilization is

$$\Delta_i = \int_{t_{i-1}}^{t_i} V(t)dt - x(t_i - t_{i-1}),$$

with respect to the expected mean use  $x$  per unit of time. Define the thresholds  $\theta_n$  ( $n = -N, \dots, N$ ) such that  $\theta_i < \theta_j$  if  $i < j$ , and  $\theta_0 = 0$ . Let also  $\theta_{-(N+1)} = -\infty$  and  $\theta_{(N+1)} = \infty$ .  $c_i$  cumulus points (positive or negative) are assigned by the ISP to the user during period  $[t_{i-1}, t_i]$  if  $\theta_{c_i} \leq \Delta_i < \theta_{c_i+1}$ .

Let  $\Lambda_n = \sum_{i=1}^n c_i$  be the sum of cumulus points assigned to the user during  $[0, t_n]$ . The ISP reacts and renegotiates the contract if  $|\Lambda_n| \geq \Theta$ .

The tariff function  $p(x)$  per unit at service level  $x$  has to be determined (the total charge is  $c(x) = xp(x)$ ). For convenience,  $p(x)$  will also be used for extra-fees: if the *observed* service level is  $x_1$ , the penalty charge is

$$\Psi(x, x_1) = c(x_1) - (c(x) + c(x_1 - x)).$$

The following requirements are inserted in order to obtain a fair scheme:

1.  $p(x) > 0$  is monotonically decreasing;  $c(x)$  is monotonically increasing.
2.  $\Psi(x, x_1) < 0$  if  $x \neq x_1$  and  $\Psi(x, x_1) = 0$  if  $x = x_1$ , so that, from the penalty charge, the user has the incentive to tell his true service level requirement.  $\Psi(x, x + \delta)$  is decreasing in  $\delta$ .
3.  $|\Psi(x, x_1)| < |\Psi(\beta x, \beta x_1)| \leq \beta |\Psi(x, x_1)|$  for  $\beta > 1$ , meaning that the penalty is higher for high bandwidths, but smaller proportionally to the expected ones.

For instance,  $p(x) = C/\sqrt{x}$  fulfills these requirements. It is suggested that no more than 3-5 thresholds are used. Moreover the number of assigned cumulus points should be "independent" of the measurement technique for determining  $x_1$ . Assuming that the stochastic process  $V(t)$  is in equilibrium, and performing  $N$  independent measurement during each interval  $[t_{i-1}, t_i]$ , a confidence interval of the  $E(V)$  can be obtained using Student distribution, at confidence level  $1 - \alpha$ , by a standard Monte Carlo method. Let  $\varepsilon_{\alpha, N}$  be the half width of the interval. Then taking  $\theta_{i+1} - \theta_i > 2\varepsilon_{\alpha, N}$  will ensure that, with probability at least  $1 - \alpha$ , the number of assigned cumulus points is not sensitive to the measurement technique.

## 4 Posted priority pricing

### 4.1 Work by Bohn et al.[5]

In their work, Bohn et al. use the 3-bit precedence field in the protocol header to introduce priorities (from 0, the lowest, to 7, the highest) in the traffic

like it was imagined (but not publicized) in the mid-80s when the NSFNET backbone was highly congested. This scheme is proposed in [5] as an interim solution before that the Internet is redesigned to incorporate protocols with bandwidth reservation, but it is worth studying.

Internet Service Providers negotiate with users some soft quotas on the total volume of traffic by specific IP Precedence levels: a quota system is introduced to discourage users from setting high precedence values in all their traffic. Another solution is to buy a total quota which is a weighted sum of the priority values in its packets per unit time. They suggest the formula

$$Q = \sum_{i=2}^6 x_i \alpha^{i-2}$$

where  $Q$  is the total quota used by the customer,  $x_i$  is the number of packets sent with priority  $i$  during the metered period and  $\alpha$  is a parameter greater than 1 (they propose  $\alpha = 2$ ). Priority levels 0 and 1 are not considered in the formula because they are free, and priority level 7 is reserved for network management.

This scheme is not directly related to pricing, but a pricing scheme can be devised by the ISP. It can also be seen as a charging scheme between the previous CPS and the next posted priority pricing.

## 4.2 Work by Cocchi et al.[8, 9]

This work uses also the 3-bit precedence field in the protocol header to introduce priorities. The model is the following. Let  $s_i$  denote a characterization of the network service received by the  $i$ th user ( $1 \leq i \leq n$ ) and  $V_i(s_i)$  denote the  $i$ th user's level of satisfaction, expressed in money, with a given network service  $s_i$  (we will give some examples later). If the user is charged an amount  $c_i$  for that service, the overall level of satisfaction is  $U_i = V_i(s_i) - c_i$ . Each user sends a request  $\sigma_i$  (not necessarily involving a call set-up). Let  $\underline{\sigma} = (\sigma_1, \dots, \sigma_n)$  and let  $s_i(\underline{\sigma})$  be the resulting network service. Define

$$\underline{\sigma}^{max} = \operatorname{argmax}_{\underline{\sigma}} \sum_{i=1}^n V_i(s_i(\underline{\sigma})) \text{ and } V_{max} = \sum_{i=1}^n V_i(s_i(\underline{\sigma}^{max}))$$

as respectively the vector maximizing the total satisfaction and the maximum total satisfaction. As each user is acting selfishly, i.e., is trying to maximize his own satisfaction  $U_i(s_i(\underline{\sigma})) = V_i(s_i(\underline{\sigma})) - c_i(\underline{\sigma})$ , the system needs to be in *Nash equilibrium*. Formally,  $\underline{\sigma}$  is a Nash equilibrium if for all  $i$  and all  $\tilde{\sigma}_i$ ,  $U_i(\underline{\sigma}) \geq U_i(\underline{\sigma}|\tilde{\sigma}_i)$  where  $(\underline{\sigma}|\tilde{\sigma}_i)$  is the vector where  $i$ th coordinate of  $\underline{\sigma}$  is replaced by  $\tilde{\sigma}_i$ . It means that user  $i$  can not increase alone his level of satisfaction. A pricing scheme is then said to be *acceptable* if  $\underline{\sigma}^{max}$  is the unique Nash equilibrium scheme. It can be easily seen that without a pricing scheme, i.e.  $c_i(\underline{\sigma}) = 0$ , the Nash equilibrium is unlikely to be met.

The scheme is then illustrated by examples. In [9], a simple two classes model is simulated on two different network topologies. The two different classes have different service priority at each switch (or node) of the network. Per-byte pricing is used with a higher price for the highest priority. The applications considered are e-mail, FTP, Telnet and Voice. The different functions  $V_i$ , following the required QoS, are

$$\begin{aligned} V_{\text{email}} &= -0.1(\text{avg. message delay (sec)}) \\ &\quad -(\% \text{ of messages not delivered in loose delay of 5 minutes}) \\ V_{\text{FTP}} &= 100(\text{average normalized throughput}) \\ V_{\text{Telnet}} &= -(\text{avg. packet round trip time (ms)})/10 \\ V_{\text{Voice}} &= -(\% \text{ of packets not obeying the tight delay of 100ms})-d/100 \end{aligned}$$

where  $d$  is the average one-way delay of voice packets (in ms). The requests  $\sigma_i$  are merely the priority settings on the packets. In the implementation, each particular application is assumed to use the same priority settings. The range of acceptable prices are given according to the topology of the network, but some exist for a wide range of network conditions.

In [8], the same kind of example is used, but in addition to the two service priority, there is a blocking priority, inducing then 4 different classes. This situation is interesting for some applications require small delays but can afford losses or inversely require no or very few losses but can afford delay. We then have four prices per byte  $p_{i,j}$ ,  $0 \leq i, j \leq 1$  where the first bit  $i$  means that the service priority flag is on or off and  $j$  gives the status of the no-drop flag.

### 4.3 Work by Honig and Steiglitz [30]

In this model,  $K$  users are assumed to compete for a resource (possibly at the gateway of a network; or directly at a switch for instance) for the same type of traffic, meaning the same type of QoS. User  $k$  wishes to send packets at rate  $\lambda_k$ , so that the total rate is  $\Lambda = \sum_{active\ k} \lambda_k$ . The QoS perception is given by a function  $D(\Lambda)$ . In [30], the delay represents the QoS, but other measures can be considered. A utility function  $u_k(\delta)$  is associated to user  $k$ , depending on the observed QoS  $\delta$ . If the price per packet is  $P$ , user  $k$  transmits his packets if and only if  $u_k(\delta) \geq P$ . In equilibrium, the QoS announced by the network must be what the user observes, that is the following fixed-point equation must be satisfied

$$D \left( \sum_{k : u_k(\delta) \geq P} \lambda_k \right) = \delta.$$

Under some assumptions ( $u_k$  monotonically decreasing and with limit 0 at  $+\infty$  and  $D$  strictly positive, finite, continuous, and monotonically increasing), it can be proved that there is a unique equilibrium for each price  $P$ . The idea is then to choose the price  $P$  maximizing the revenue  $R = P\Lambda$ . Some examples are provided.

As extensions, multiple priorities and time of day pricing are discussed.

### 4.4 Work by Marbach [48, 49, 50]

This work is devoted to DiffServ, where packet classes are served according to a given priority. Prices per sent-packet are static. Indeed, it is argued that, by charging for all *submitted* packets, the users receive the incentive to reduce their rates during periods of congestion, as they pay for lost packets.

The mathematical model is considering a single link and is the following. The time is discretized, divided into slots. During each slot, the link has the ability to serve  $C$  packets. It is assumed that packets not served during the slot are lost. They are  $N$  different (and ordered) priority classes, where 1 is the lowest priority.  $R$  users are supposed to compete for the link access. Let  $u_i$  be the price charged for a class- $i$  packet submitted for access (of course  $u_i < u_j$  if  $i < j$ ) and  $d_r(i)$  be the number of class- $i$  packets that user  $r$  is submitting in a given time slot. User  $r$ 's whole allocation is given by the

vector  $d_r = (d_r(1), \dots, d_r(N))$ . The number of submitted class- $i$  packets is  $d(i) = \sum_{r=1}^R d_r(i)$  and  $d = (d(1), \dots, d(N))$  is the aggregated allocation.

Let  $i^*$  be the priority class such that  $\sum_{i=i^*+1}^N d(i) < C$  and  $\sum_{i=i^*}^N d(i) \geq C$ . Packets with priority  $i > i^*$  are served (say, with probability  $P_{tr}(i, d) = 1$ ), those with priority  $i < i^*$  are lost (say, with probability  $P_{tr}(i, d) = 0$ ) and those of class  $i^*$  are served with probability

$$P_{tr}(i^*, d) = \frac{(C - \sum_{i=i^*+1}^N d(i))}{d(i^*)}.$$

User  $r$  throughput is then

$$x_r = \sum_{i=1}^N d_r(i) P_{tr}(i, d).$$

A utility function  $U_r(x_r)$  is associated to user  $r$ .  $U_r$  is assumed to be increasing, bounded, strictly concave and twice differentiable. The users are assumed to play a non-cooperative game, where user  $r$  chooses allocation  $d_r^*$  such that

$$d_r^* = \operatorname{argmax}_{d_r} \left( U_r(x_r) - \sum_{i=1}^N d_r(i) u_i \right).$$

In equilibrium, this happens for all users. If we suppose without loss of generality that the total demand at price  $u_1$ ,  $D(u_1)$ , exceeds  $C$  and that  $D(u_i) > 0 \forall i$ , then

- there exists an equilibrium. If there is a class  $i_0$  such that  $D(u_{i_0}) > C > D(u_{i_0+1})$ , the equilibrium is unique.
- $d_r^*(i) = 0 \forall i \notin \{i_0, i_0+1\}$ ;  $P_{tr}(i_0, d^*) \geq u_{i_0}/u_{i_0+1}$  where  $u_{N+1} = \max_r U_r'(0)$ ; and  $x_r^* = D_r(u^*)$  with  $u^* = u_{i_0}/P_{tr}(i_0, d^*)$ .

In [48], the game is played dynamically. A gradient algorithm is used to prevent from oscillations. In [49], the model is extended to bursty traffic.

## 5 Non-posted Priority pricing

### 5.1 Work by Mendelson and Whang [55], and by Ha [27]

The work described here was not dedicated to the Internet management. But, even if some points are irrelevant for our concern, it is worth studying.

Mendelson and Whang consider a pricing scheme for a multi-class M/M/1 queue (which can be easily extended to a M/G/1 queue if all job classes have the same coefficient of variation). Arrivals of class- $i$  jobs ( $1 \leq i \leq R$ ) to the system reflect the aggregation of infinitesimal users' job flow. The arrival rate is  $\lambda_i$ . The value function of class- $i$  jobs  $V_i(\lambda_i)$ , representing the gross value gained by class- $i$  users per unit of time, is assumed to be differentiable, nondecreasing and concave on  $\lambda_i$ .  $\lambda_i$  and the "full price"  $z$  are related in the following way:  $\lambda_i = D_i(z) = (1 - F_i(z))\Lambda_i$  where  $\Lambda_i$  is the maximum potential arrival rate of class  $i$  and  $F_i(\cdot)$  is the distribution function of the service valuation. Inverting this function, we have  $V_i'(\lambda_i) = D_i^{-1}(\lambda_i)$ . Let  $\underline{\lambda} = (\lambda_1, \dots, \lambda_R)$ . The total expected value function is

$$V(\underline{\lambda}) = \sum_{i=1}^R V_i(\lambda_i).$$

Each class- $i$  job is characterized by a delay cost of  $v_i$  per unit time. Class- $i$  jobs are assumed to be served following an exponential distribution with mean  $c_i$  and the priority policy of the server is supposed to be nonpreemptive. Assume also, without loss of generality, that the classes are ordered from the highest to the lowest priority so that the expected average delay cost per unit time is minimized, i.e.,

$$\frac{v_1}{c_1} \geq \frac{v_2}{c_2} \geq \dots \geq \frac{v_R}{c_R}.$$

The idea is to maximize the expected net value of the jobs processed by the system, i.e. to find

$$\max_{\underline{\lambda}} \left\{ V(\underline{\lambda}) - \sum_{i=1}^R v_i L_i(\underline{\lambda}) \right\} \quad (1)$$

where  $L_i$  is the mean number of class- $i$  jobs in the system in steady-state. The administrator set the price vector

$$\underline{p} = (p_1, \dots, p_R)$$



where  $p_i$  is the price charged to a class- $i$  job. If class- $i$  demand relationship (setting that at equilibrium the marginal value will be indifferent between joining and not joining the system) is  $V_i'(\lambda) = p_i + v_i W_i(\lambda)$ , it is proved in [55] that the optimal price per class- $i$  job is given by

$$p_i^* = \sum_{j=1}^R v_j \lambda_j^* \frac{\mathcal{D}W_j(\underline{\lambda}^*)}{\mathcal{D}\lambda_i}$$

where  $W_j$  is the expected delay of a class- $j$  job and where  $\underline{\lambda}^*$  maximizes (1).

In the homogeneous case, i.e.,  $c_1 = \dots = c_R$ , we have  $L_i(\underline{\lambda}) = \frac{\lambda_i S_R}{S_{i-1} S_i} + \lambda_i$  and  $W_i(\underline{\lambda}) = \frac{S_R}{S_{i-1} S_i} + 1$  where  $S_0 = 0$ ,  $S_i = \sum_{j=1}^i \lambda_j$  and  $\bar{S}_i = 1 - S_i$ . We can then get explicitly the optimal prices:

$$p_i^* = \sum_{k=1}^R \frac{\lambda_k^* v_k}{\bar{S}_k \bar{S}_{k-1}} + \sum_{k=1}^R \frac{\lambda_k^* v_k W_k^q + \lambda_{k+1}^* v_{k+1} W_{k+1}^q}{\bar{S}_k}$$

where  $W_k^q = W_k - 1$  is the expected waiting time of a class- $k$  job in the queue and  $\lambda_{R+1}^* = v_{R+1} = W_{R+1}^q = 0$ .

The problem here is that the prices are determined on a centralized basis, which is practically irrelevant. More specifically, both user and system administrator know  $(v_i, V_i, c_i) \forall i$ , but only users know their *real* class membership. To cope with this problem, Mendelson and Whang consider priority-dependent pricing schemes. The idea is to obtain a *Nash equilibrium*, i.e. that no user, by unilaterally changing its own request, can increase his own net value. This property is decomposed in *incentive-compatibility* which means that it is in all users' interest to classify their jobs in their correct priority class, and *optimality* which means that the resulting arrival rates maximize the expected net value of the system as a whole. Optimality and incentive-compatibility are obtained when using the optimal prices in the homogeneous case and when a class- $i$  user decides not to enter the system if

$$\min_{1 \leq j \leq R} \{0, p_j + v_i W_j(\underline{\lambda}(\underline{p})) - V_i'(\lambda_i)\} = 0$$

and to join the system otherwise.

Unfortunately, the incentive-compatibility is not valid anymore in the heterogeneous case. The previous posted charging mechanism should take into

account additional information such as the actual processing time of the job. We then have a priority and time-dependent pricing scheme. If we have

$$p_i(t) = A_i t + (1/2) B t^2$$

with

$$B = \sum_{k=1}^R \frac{v_k \lambda_k^*}{\bar{S}_{k-1} \bar{S}_k}$$

and

$$A_i = \frac{a_i}{\bar{S}_{i-1} \bar{S}_i^2} + \sum_{k=i+1}^R a_k \left( \frac{1}{\bar{S}_{k-1} \bar{S}_k} + \frac{1}{\bar{S}_{k-1} \bar{S}_k^2} \right)$$

where  $a_i = v_i \lambda_i^* \sum_{k=1}^i c_k^2 \lambda_k^*$ , then the pricing scheme is optimal and incentive-compatible. Note that it is decomposed in a basic charge (corresponding to the lowest priority charge) and a priority surcharge (proportional to the processing time).

In [27], A.Y. Ha extends the previous work to the case when service requirements are controllable by the customers. Then, each customer decides whether to request service from the facility and, if desirable, determines his service requirement. The case of the  $M/G/s$  processor sharing queue is investigated, for which the optimal prices are found to be two-parts linear in time in the system. The first-come-first served  $M/G/1$  queue is also studied and a quadratic price is also obtained.

## 5.2 Work by Gupta et al.

In [25, 26], Gupta, Stahl and Whinston develop also a priority pricing scheme. They first argue that the posted priority pricing scheme of Bohn et al. may be a lack of incentive to provide multiple precedence networks (i.e., the providers may not be appropriately rewarded) and that we must look at the context in which the applications are used, not just to categorize them.

In [25], a four priority classes model is introduced, where the highest priority is for real-time services with no tolerance for lost packets, the second class is for real-time services that are relatively tolerant to lost packets and the two lowest priority classes for two levels of best effort service (to provide a finer division of delay requirements). In [26], the number of classes is kept general.

The price at a particular server for a particular class is represented by the following system of equations:

$$r_{mk}(q) = \sum_l [\mathcal{D}\Omega_l / \mathcal{D}\chi_{mkq}] \sum_i \sum_j \delta_{ij} x_{ijlm} \quad (2)$$

where

- $r_{mk}(q)$  is the price of a job sized  $q$  at server  $m$  for priority class  $k$
- $\chi_{mkq}$  is the arrival rate of job sized  $q$  at server  $m$  in priority class  $k$
- $\Omega_l$  is a continuously differentiable, strictly increasing function of arrival rate  $\chi_{mkq}$  and capacity  $v_m$  which provides the waiting time at a server  $m$  for priority class  $l$
- $\delta_{ij}$  is the delay cost parameter of consumer  $i$  for service  $j$
- $x_{ijlm}$  is the flow rate of service  $j$  for consumer  $i$  with priority  $k$  at server  $m$ .

$[\mathcal{D}\Omega_l / \mathcal{D}\chi_{mkq}]$  is the derivative of the waiting time and  $\sum_i \sum_j \delta_{ij} x_{ijlm}$  is the accumulated delay cost of the system. This kind of priority pricing prevents that the “highest priority can preempt all the available capacity” (as said in [7]) in the case of posted priority of previous subsections.

In [26], a general mathematical model is introduced and it is shown that this choice maximizes a system-wide welfare stochastic allocation function.

The prices are computed using the following iterative equation:

$$r_{mk}^{t+1} = \alpha \hat{r}_{mk}^{t+1} + (1 - \alpha) r_{mk}^t$$

where

- $\alpha$  is a real number between 0 and 1. The authors suggest to take  $\alpha = 0.1$ .
- $\hat{r}_{mk}^{t+1}$  is the estimated new price at time  $t + 1$  using Equation (2)
- $r_{mk}^t$  is the implemented price during the time interval  $(t, t + 1)$ .

Many experiments are performed using a simulation platform.

## 6 Smart market: auction in the network

### 6.1 Smart market of McKie-Mason and Varian [53, 54]

In their paper about internet history, cost and pricing [53], McKie-Mason and Varian argue that posted priority pricing of Section 4 is not a good solution. Indeed, if the network is at capacity, some users with high willingness-to-pay may be unable to access the network. Pricing by time of day attempts to achieve this goal but does not allocate efficiently the available bandwidth.

McKie-Mason and Varian suggest the use of a “smart market”, which is actually a variation of the Vickrey auction. Each packet is given a bid representing the user’s willingness to pay. The packets are given a priority at each node of the network according to this bid. Using Vickrey auction, if the network is uncongested, the price is zero whereas if there is congestion, the charge is the willingness to pay of the lowest priority packet admitted.

Unfortunately, smart market is not an ideal solution. As said in [53], current TCP/IP would not support a smart market. Moreover it requires complicated systems to conduct auctions for individual packets. The model was more an incentive for further research than a solution.

### 6.2 Progressive Second Price (PSP) Auction

In [43, 68, 69, 70, 75], the costly auctions for individual packets are replaced by auctions for bandwidth during intervals of time. A good analysis of this scheme based on game theory is provided, including fairness properties. As stated in [68], "in market-based approaches, no precise model need be assumed [...], the seller does not require a priori demand information". The behavior of the system is then essentially real-time, and not model-based.

To briefly explain how the auction works, consider a single resource of capacity  $Q$  and  $I$  players competing for it. Player  $i$ 's bid is  $s_i = (q_i, p_i)$  where  $q_i$  is the capacity the player  $i$  is looking for and  $p_i$  is the unit price he is proposing. A bid profile is  $s = (s_1, \dots, s_I)$ . Let  $s_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_I)$  be the profile where player  $i$ 's bid is excluded from the game. For  $y \geq 0$  define

$$\underline{Q}_i(y; s_{-i}) = \left[ Q - \sum_{p_k \geq y, k \neq i} q_k \right]^+.$$

The progressive second price allocation rule gives to player  $i$  a bandwidth

$$a_i(s) = \min(q_i, \underline{Q}_i(p_i; s_{-i}))$$

and set the total cost to

$$c_i(s) = \sum_{j \neq i} p_j [a_j(0; s_{-i}) - a_j(s_i; s_{-i})].$$

Thus the highest bids are allocated the desired quantity and the cost is given by the declared willingness to pay (bids) of the users who are excluded by  $i$ 's presence.

Assume that player  $i$  attempts to maximize his utility  $u_i(s) = \theta_i(a_i(s)) - c_i(s)$  where  $\theta_i$  is the valuation function that player  $i$  gives to his allocation. Under some smoothness assumptions on  $\theta_i$  and having a bid fee  $\varepsilon$  each time a player submits a bid, it is stated that if for all  $i$  player  $i$  bids  $(v_i, \omega_i = \theta'_i(v_i))$  with

$$v_i = \left[ \sup \left\{ z : z \leq Q_i(\theta'_i(v_i), s_{-i}) \text{ and } \sum_{j \neq i} p_j [a_j(0; s_{-i}) - a_j(z; s_{-i})] \leq b_i \right\} - \varepsilon / \theta'_i(0) \right]^+$$

where  $Q_i(y; s_{-i}) = [Q - \sum_{p_k > y, k \neq i} q_k]^+$  and  $b_i$  is the budget constraint, then convergence, efficiency and fairness issues are solved (the property in [75] that the equal-bid case (when the total required bandwidth at this unit price is not available) does not occur in the PSP scheme).

The game is extended in [68, 70] to networked auctions and the same properties are obtained. In this networked game, players can be raw bandwidth sellers, end-users or service providers buying and selling bandwidth to each others. Each player is acting like in the single node case, trying to optimize its utility  $u_i = \theta_i \circ e_i(a) - \sum_j c_i^j$  where  $e_i$  is a function called the *expected bottleneck* depending on the type of player and  $c_i^j$  is the total cost charged to player  $i$  by seller  $j$ .

In [10], simultaneous multi-unit descending-price auctions (or Dutch auctions) with different decreasing speed are used. Indeed, they argue that, among other drawbacks, in PSP auction each player splits equally his bid among links, which might not be right (depending on the congestion levels). The mechanism allows then that each user buys the same quantity of bandwidth capacity

in all the links. According to experimental results, social welfare is improved with respect to PSP.

In [65], two new auction schemes are designed: delta auction which allows bids to take place continuously in order to prevent from the additive setup delay (at each node) and Connection-Holder-is-Preferred-Scheme (CHiPS), based on RSVP protocol, for which holders of already running connections are preferred and are given a second chance if their actual bid is exceeded by a new one.

## 7 Expected capacity [7]

In [7], Clark is also discussing how to charge the Internet. Like many authors, he is more setting open issues than solving them. One of his questions about priority pricing is the following: “the effect of priority queuing is to build up a queue of lower-priority packets which will cause packets in this class to be preferentially dropped due to queue overflow. While dropped packets will be retransmitted, the rate adaptation of TCP translates these losses into a reduction in sending packets for these flows of packets”. Moreover he says that there is no obvious way to relate a particular priority with a particular achieved service. He introduces then his notion of *expected capacity*. The mechanism works as follows. At the network access, packets are flagged (*in* or *out*) depending on the fact that the incoming stream is inside or outside of the profile of the expected capacity (without any traffic shaping). When there is a point of congestion, *out*-tagged packets receive a congestion pushback notification (dropping or explicit congestion notification (ECN)). During periods of congestion, each sender executes a TCP algorithm which receives a congestion indication when it exceeds its expected capacity and starts to send packets that are flagged *out*.

As said in [7], this scheme can also be implemented in a heterogeneous network of multi-provider Internets where cooperating groups of providers make contracts to carry each other’s traffic; when too many packets are marked in according to the contract, they can be shifted to out or they can be charged according to some formula. Some dynamic tagging can also be implemented as done with smart markets of McKie-Mason and Varian.

Unfortunately, some problems need to be solved to implement efficiently this scheme. First, depending on the application, the customer can be the sender or the receiver. The scheme previously described is working if the customer is the sender. If he is the receiver, there is a need to design a complex protocol by which the sender is informed of the expected capacity contract, which can be also quite complex (to keep flexibility of contracts). Second, what about multicast when each receiver has a different expected capacity?

## 8 Charging for elastic traffic based on transfer rate [3, 33, 36, 46, 45]

### 8.1 Work by Kelly et al.

The model presented here allows to combine different elastic traffics [32, 33, 36] where the rates are proportional to the willingness to pay of each user. The model is the following. Consider a set of  $J$  resources with a capacity of  $C_j$  for resource  $j$ . A route  $r$  is a non-empty subset of  $J$  and  $R$  is the set of possible routes. Let  $A_{jr} = 1$  if  $j \in r$  and 0 otherwise and define  $A$  as  $A = (A_{jr})$ . If each route is associated to a user  $r$ , let  $U_r(x_r)$  be the utility function of the user when the flow rate is  $x_r$  for user  $r$ .  $U_r$  is assumed to be an increasing, strictly concave and continuously differentiable function. Let  $U = (U_r(\cdot), r \in R)$  and  $C = (C_j, j \in J)$ . From the system point of view, the idea is to maximize

$$\sum_{r \in R} U_r(x_r) \quad (3)$$

subject to  $Ax \leq C$  and  $x \geq 0$ . From the user point of view, the idea is to maximize

$$U_r\left(\frac{\omega_r}{\lambda_r}\right) - \omega_r \quad (4)$$

over  $\omega_r \geq 0$ ; here the flow rate is  $x_r = \omega_r \lambda_r$  where  $\omega_r$  is the amount that user  $r$  is willing to pay per unit time and  $\lambda_r$  is the charge per unit flow and unit time for user  $r$ . Assume that the network knows  $\omega = (\omega_r, r \in R)$  and attempts to maximize

$$\sum_{r \in R} \omega_r \log x_r \quad (5)$$

subject to  $Ax \leq C$  and  $x \geq 0$ . This last assumption is very convenient because it allows to compute optimal flow rates very easily. Indeed, it is shown in [32, 33, 36] that there always exist vectors  $\lambda$ ,  $\omega$  and  $x$  satisfying  $\omega_r = x_r/\lambda_r \forall r \in R$  such that  $\omega_r$  maximizes (4),  $x$  maximizes (5) and then  $x$  is the unique solution maximizing (3).

It is also shown that the vector of rates  $x$  per unit charge is *proportionally fair*, that is, if  $x \geq 0$  and  $Ax \leq C$ , and for any other feasible vector  $x^*$ , the aggregate proportional change is zero or negative

$$\sum_{r \in R} \omega_r \frac{x_r^* - x_r}{x_r} \leq 0.$$

Even if solving this problem is mathematically tractable, the maximization of (5) needs to be done on a centralized basis, which is undesirable. In the following is explained how to proceed on a decentralized basis. Consider the system of differential equations

$$\frac{d}{dt} x_r(t) = \kappa_r \left( \omega_r(t) - x_r(t) \sum_{j \in r} \mu_j(t) \right) \quad (6)$$

where

$$\mu_j(t) = p_j \left( \sum_{s: j \in s} x_s(t) \right)$$

is the shadow price per unit flow through  $j$  and  $p_j(t)$  is the derivative of the rate at which cost is incurred at resource  $j$  when the load through it is  $y$ . The motivation behind these equations is the following. If resource  $j$  generates a continuous stream of feedback signal at rate  $yp_j(y)$  when the total flow through resource  $j$  is  $y$ ; that resource  $j$  sends a proportion  $x_r/y$  of these feedback signals to a user  $r$  with a flow of rate  $x_r$  through resource  $j$ ; and that user  $r$  views each feedback signal as a congestion indication requiring some reduction of flow  $x_r$ . It is then a flow-control algorithm. It is shown using Lyapunov functions that the system of differential equations has a unique value  $x$  such that  $x_r = \omega_r / \sum_{j \in r} \mu_j$  arbitrarily closely approximates the optimization of problem (5). Some stochastic perturbations of equation (6) are also analysed in [36].



Equation (6) shares several characteristics with TCP but presents also several differences as pointed out in [33]. In TCP, congestion indication is from dropped or marked packets. There is then here two multiplicative effects. Anyway, it is shown that multiple TCP can be modeled by the system of differential equations

$$\frac{d}{dt}x_r(t) = \frac{m_r}{T_r^2} \left( \frac{m_r}{T_r^2} + \frac{x_r(t)^2}{2m_r} \right) \sum_{j \in r} \mu_j(t) \quad (7)$$

and can be viewed as acting as if the utility function of user  $r$  is

$$\frac{\sqrt{2}m_r}{T_r} \arctan \left( \frac{x_r T_r}{\sqrt{2}m_r} \right)$$

where  $T_r$  is the round trip time for the connection of user  $r$  and  $m_r$  is a parameter which would *inter alia* multiply by  $m$  the rate of additive increase and make  $1 - 1/2m$  the multiplicative decrease factor in Jacobson's TCP algorithm. The stable point is then such that  $\forall r$

$$x_r = \frac{m_r}{T_r} \left( \frac{2(1 - p_r)}{p_r} \right)^{1/2}$$

where  $p_r = \sum_{j \in r} p_j$ . Note that this conclusion cannot be reached when users or the network have routing choices.

Each customer can use intelligent agents [11] in order to optimize his willingness to pay according to the network congestion status.

In [39, 40], the necessary feedback to the users who adjust their rates is based on window-based congestion control, which is practically easy to do when connections are using TCP. The method is proved to give optimal values. It is shown that the solution solves the same problem than the one of Kelly et al.. Other implementations of the scheme are presented in [16, 38, 37, 41], giving some scenari and algorithms for user adaptation and network feedback signals for flow control.

## 8.2 Work by Low et al.

In [3, 46, 45] Low et al. study the same kind of problem than Kelly et al. (used mainly for ABR in ATM networks rather than for TCP in the Internet)

and they get very similar solutions. The main difference is that in Low et al.'s work, the users decide their rates and pay for it whereas in Kelly et al.'s work, the users decide their payment and receive what the network allocates. Here also, they use a decentralized algorithm to set up prices adapted to changing network conditions. As in previous subsection, we have a set  $L$  of unidirectional links of capacities  $c_l$ ,  $l \in L$ , a set  $S$  of sources characterized by utility function  $U_s(x_s)$  concave, increasing in its transmission rate  $x_s$ . The system is willing to maximize

$$\sum_{s \in S} U_s(x_s)$$

over  $x_s$  subject to capacity constraints. The problem is also decomposed and the following synchronous algorithm is used in [46]:

1. Each link receives the rates  $x_s(t)$  if  $s$ 's route is through link  $l$ .
2. Each link  $l$  calculates its price  $p_l(t+1)$  for a unit of bandwidth (in order to optimize the benefits) using the gradient projection algorithm

$$p_l(t+1) = [p_l(t) + \gamma(x^l(t) - c_l)]^+ \quad (8)$$

where  $\gamma$  is a stepsize.

3. Each link communicates  $p_l(t+1)$  to each source whose route is through link  $l$ .

Then the algorithm for each source is:

1. Each source is fed back the price  $p^s = \sum_{L(s)} p_l$  where  $L(s)$  is the set of links that  $s$  uses.
2. The source chooses then its transmission rate  $x_s$  (in an interval  $(m_s, M_s]$ ) which maximizes its benefit

$$U_s(x_s) - p^s(t)x_s.$$

3. These rates  $x_s(t+1)$  are send to the links which calculate again new prices and so on.

The algorithm approaches a price vector  $(p_l^*, l \in L)$  that aligns individual and system optimality with fairness properties. In [3], the gradient projection method is replaced by the Newton method which typically converges much faster. Equation (8) is then replaced by

$$p_l(t+1) = [p_l(t) + \gamma H_l^{-1}(t)(x^l(t) - c_l)]^+ \quad (9)$$

where  $H$  is a Hessian matrix (see [3] for details). In [2], the equation is replaced by

$$p_l(t+1) = [p_l(t) + \gamma(\alpha_l b_l(t) + x^l(t) - c_l)]^+ \quad (10)$$

where  $\alpha_l$  is a constant and  $b_l(t)$  is the buffer backlog at link  $l$ .

The model is extended to the asynchronous case where the updates at the sources and the links are not synchronized, which better resembles the reality of large networks. The communication between sources and links is also greatly simplified as follows. In [45], the links estimate source rates using local information without leaving the optimality property. In [42, 2], the communication from links to sources is accomplished using the proposed ECN (Explicit Congestion Notification) bit in the IP header. These modifications lead to a flow control scheme called REM (Random Early Marking), a variant of RED, and a stochastic version of the previous algorithm: link  $l$  marks an arriving packet with probability  $m_l(t) = 1 - \Phi^{-p_l(t)}$  (with  $\Phi > 1$ ). This leads to  $m^s(t) = 1 - \Phi^{-p^s(t)}$ . Inverting this equation,  $p^s(t)$  is estimated by  $\hat{p}^s(t) = -\log_\Phi(1 - \hat{m}^s(t))$  where  $\hat{m}^s(t)$  is the fraction of marked packets (known by usual acknowledgement). Stability, performance and robustness of this version of the algorithm is studied in [60] using a continuous time version of the dynamics.

## 9 Conclusions

In this paper we have surveyed usage-based pricing schemes without bandwidth reservation. All these schemes have their own advantages, ranging from implementation simplicity to fairness issue. An interesting problem would be to compare (mathematically and in practice) their respective costs and benefits on a simple network in order to point out the one which is likely to perform the best.

Also, other issues worth to be studied. First, an interesting area of research is the pricing of Weighted Fair Queueing schemes. According to Clark [7], this mechanism would only achieve local equality inside one switch. For example in the multicast case, what does congestion along one path as to do with congestion along another? This also shows that the multicast case [4, 22, 23] needs more attention, like in [24], where an adaptation of pricing based on transfer rates is applied to multicast flows and fairness properties are obtained. Next, as pointed out in [72], optimality paradigm is not a panacea and more attention needs to be given to architecture and structure.

## References

- [1] L. Anania and R.J. Solomon. Flat- The Minimalist Price. In Lee W. McKnight and Joseph P. Bailey, editors, *Internet Economics*, pages 91–118. MIT Press, 1997.
- [2] S. Athuraliya and S.H. Low. Optimization Flow Control, II: Implementation. Technical report, 2000.
- [3] S. Athuraliya and S.H. Low. Optimization Flow Control with Newton-Like Algorithm. *Telecommunication Systems*, 13, 2000.
- [4] A. Basu and S.J. Golestani. Estimation of Receiver Round Trip Times in Multicast Communications. Technical report, Bell Laboratories. <http://www.bell-labs.com/user/golestani/rtt.ps>.
- [5] R. Bohn, H.W. Braun, K.C. Claffy, and S. Wolff. Mitigating the Coming Internet Crunch: Multiple service levels via Precedence. Technical report, University of California - San Diego, 1993.
- [6] T. Bonald and L. Massoulié. Impact of Fairness on Internet Performance. In *Proceedings of ACM Sigmetrics 2001*, 2001.
- [7] D.D. Clark. Internet Cost Allocation and Pricing. In Lee W. McKnight and Joseph P. Bailey, editors, *Internet Economics*, pages 215–252. MIT Press, 1997.

- 
- [8] R. Cocchi, D. Estrin, S. Shenker, and L. Zhang. A Study of Priority Pricing in Multiple Service Class Networks. In *Proceedings of SIGCOMM'91*, pages 123–130, 1991.
  - [9] R. Cocchi, D. Estrin, S. Shenker, and L. Zhang. Pricing in Computer Networks: Motivation, Formulation and Example. *IEEE/ACM Transactions on Networking*, 1(6):614–627, 1993.
  - [10] C. Courcoubetis, M.P. Dramitinos, and G.D. Stamoulis. An auction mechanism for bandwidth allocation over paths. Technical report, Athens University of Economics and Business, 2001.
  - [11] C. Courcoubetis, G.D. Stamoulis, C. Manolakis, and F.P. Kelly. An intelligent agent for optimizing QoS-for-money in priced ABR connections. *Telecommunications Systems*, 2000.
  - [12] L.A. DaSilva. Pricing of QoS-Enabled Networks: A Survey. *IEEE Communications Surveys & Tutorials*, 3(2), 2000.
  - [13] P. Dolan. Internet Pricing. is the end of the World Wide Wait in view? *Communications & Strategies*, 37:15–46, 2000.
  - [14] M. Falkner, M. Devetsikiotis, and I. Lambadaris. An Overview of Pricing Concepts for Broadband IP Networks. *IEEE Communications Surveys & Tutorials*, 3(2), 2000.
  - [15] Z. Fan. Pricing and provisioning for guaranteed internet services. In P. Lorenz, editor, *ICN 2001*, volume 2093 of *Lecture Notes in Computer Science*, pages 55–64. Springer-Verlag, 2001.
  - [16] A. Ganesh, K. Laevens, and R. Steinberg. Dynamics of congestion pricing. Technical Report 70, Microsoft Research Limited, Cambridge, UK, 2000.
  - [17] R. Gibbens, R. Mason, and R. Steinberg. Internet service classes under competition. *IEEE Journal on Selected Areas in Communications*, 18(12):2490–2498, 2000.

- 
- [18] R.J. Gibbens and F.P. Kelly. Measurement-based connection admission control. In *Proceedings of the 15th International Teletraffic Congress*, 1997.
  - [19] R.J. Gibbens and F.P. Kelly. Distributed connection acceptance control for a connectionless network. In *Proceedings of the 16th International Teletraffic Congress*, 1999.
  - [20] R.J. Gibbens and F.P. Kelly. resource pricing and the evolution of congestion control. *Automatica*, 35:1969–1985, 1999.
  - [21] R.J. Gibbens, S.K. Sargood, F.P. Kelly, H. Azmoodeh, R. Macfadyen, and N. Macfadyen. An Approach to Service Level Agreements for IP networks with Differential Services.
  - [22] S.J. Golestani and S. Bhattacharyya. A Class of End-to-End Congestion Control Algorithms for the Internet. In *Proceedings of ICNP 98*, October 1998.
  - [23] S.J. Golestani and K.K. Sabnani. Fundamental Observations on Multicast Congestion Control in the Internet. In *Proceedings of IEEE INFOCOM 99*, March 1999.
  - [24] E.E. Graves, R. Srikant, and D. Towsley. Decentralized Computation of Weighted Max-Min Fair Bandwidth Allocation in Networks with Multicast Flows. In S. Palazzo, editor, *IWDC 2001*, volume 2170 of *Lecture Notes in Computer Science*, pages 326–342. Springer-Verlag, 2001.
  - [25] A. Gupta, D.O. Stahl, and A.B. Whinston. Priority Pricing of Integrated Services Networks. In Lee W. McKnight and Joseph P. Bailey, editors, *Internet Economics*, pages 323–352. MIT Press, 1997.
  - [26] A. Gupta, D.O. Stahl, and A.B. Whinston. A stochastic equilibrium model of Internet pricing. *Journal of Economic Dynamics and Control*, 21:697–722, 1997.
  - [27] A.Y. Ha. Optimal pricing that coordinates queues with customer-chosen service requirements. *Management Science*, 47(7):915–930, 2001.

- 
- [28] D. Hazlett. An interim economic solution to internet congestion. *Social Science Computer Review*, 15(2):181–189, 1997.
- [29] T. Henderson, J. Crowcroft, and S. Bhatti. Congestion Pricing. Paying Your Way in Communication Networks. *IEEE Internet Computing*, September/October:85–89, 2001.
- [30] M.L. Honig and K. Steiglitz. Usage-based pricing of packet data generated by a heterogeneous user population. In *Proceedings of IEEE INFOCOM 95*, pages 867–874, 1995.
- [31] F.P. Kelly. Note on effective bandwidths. In F.P. Kelly, S. Zachary, and I.B. Ziedins, editors, *Stochastic Networks: Theory and Applications*, volume 4 of *Royal Statistical Society Lecture Notes Series*, pages 141–168. Oxford University Press, 1996.
- [32] F.P. Kelly. Charging and rate control for elastic traffic. *European transactions on Telecommunications*, 8:33–37, 1997.
- [33] F.P. Kelly. Mathematical modelling of the Internet. In *Proceedings of the Fourth International Congress on Industrial and Applied Mathematics*, 2000.
- [34] F.P. Kelly. Models for a self-managed Internet. *Philosophical Transactions of the Royal Society*, A358, 2000.
- [35] F.P. Kelly, P.B. Key, and S. Zachary. Distributed Acceptance Control. *IEEE Journal on Selected Areas in Communications*, 18, 2000.
- [36] F.P. Kelly, A.K. Mauloo, and D.K.H. Tan. Rate control in communication networks: shadow prices, proportional fairness and stability. *Journal of the Operational Research Society*, 49:237–252, 1998.
- [37] P. Key and L. Massoulié. User policies in a network implementing congestion pricing. Technical report, Microsoft Research Limited, Cambridge, UK, 1999.
- [38] P. Key and D.R. McAuley. Differential QoS and Pricing in Networks: where flow-control meets game theory. In *IEE Proceedings*, 1999.

- 
- [39] R.J. La and V. Anantharam. Charge-Sensitive TCP and Rate Control in the Internet. In *Proceedings of IEEE INFOCOM 2000*, 2000.
  - [40] R.J. La and V. Anantharam. Window-Based Control with Heterogeneous Users. In *Proceedings of IEEE INFOCOM 2001*, 2001.
  - [41] K. Laevens, P. Key, and D. McAuley. An ecn-based end-to-end congestion-control framework: experiments and evaluation. Technical Report 104, Microsoft Research Limited, Cambridge, UK, 2000.
  - [42] D.E. Lapsley and S.H. Low. Random early marking: An optimisation approach to internet congestion control. In *Proceedings of IEEE ICON'99*, 1999.
  - [43] A.A. Lazar and N. Semret. Design and Analysis of the Progressive Second Price Auction for Network Bandwidth Sharing. *To appear in Telecommunication Systems*, 13, 2001. <http://comet.columbia.edu/~nemo/telecomsys.pdf>.
  - [44] X. Lin and N.B. Shroff. Pricing-based control of large networks. In S. Palazzo, editor, *IWDC 2001*, volume 2170 of *Lecture Notes in Computer Science*, pages 212–231. Springer-Verlag, 2001.
  - [45] S.H. Low. Optimization Flow Control with On-line Measurement or Multiple Paths. In *Proceedings of the 16th International Teletraffic Congress*, 1999.
  - [46] S.H. Low and D.E. Lapsley. Optimization Flow Control, I: Basic Algorithm and Convergence. *IEEE/ACM Transactions on Networking*, 7(6), 1999.
  - [47] S.H. Low, F. Paganini, and J.C. Doyle. Internet Congestion Control. *IEEE Control Systems Magazine*, 2002.
  - [48] P. Marbach. Differentiated Services Networks: Pricing and Software Agents. Technical Report CSRG-422, Department of Computer Science, University of Toronto, 2001.



- 
- [49] P. Marbach. Pricing Differentiated Services Networks: Bursty Traffic. In *Proceedings of IEEE INFOCOM 2001*, 2001.
- [50] P. Marbach. The Role of Pricing in Differentiated Services Networks. Technical Report CSRG-421, Department of Computer Science, University of Toronto, 2001.
- [51] L. Massoulié and J. Roberts. Arguments in favour of admission control for TCP flows. In *Proceedings of the 16th International Teletraffic Congress*, 1999.
- [52] L. Massoulié and J. Roberts. Bandwidth sharing: objectives and algorithms. In *Proceedings of IEEE INFOCOM 99*, 1999.
- [53] J.K. McKie-Mason and H.R. Varian. Some Economics of the Internet. Technical report, University of Michigan, November 1993. <http://wueconb.wustl.edu:8089/eps/comp/papers:9401/9401001.pdf>.
- [54] J.K. McKie-Mason and H.R. Varian. Pricing Congestible Network Resources. *IEEE Journal on Selected Areas in Communications*, 13:1141–1149, 1995.
- [55] H. Mendelson and S. Whang. Optimal incentive-compatible priority pricing for the M/M/1 queue. *Operations Research*, 38(5):870–883, 1990.
- [56] D. Mitra, K.G. Ramakrishnan, and Q. Wang. Combined economic modeling and traffic engineering: Joint optimization of pricing and routing in multi-service networks. In *Proceedings of the 17th International Teletraffic Congress*, 2001.
- [57] J. Mo and J. Walrand. Fair end-to-end window-based congestion control. In *Proceedings of SPIE'98*, October 1998.
- [58] A. Odlyzko. Paris Metro Pricing for the Internet. In *ACM Conference on Electronic Commerce (EC'99)*, pages 140–147, 1999.
- [59] A.M. Odlyzko. The current state and likely evolution of the Internet. In *Proceedings of Globecom'99*, pages 1869–1875, 1999.

- 
- [60] F. Paganini. Flow control via pricing: a feedback perspective. In *Proceedings of the 2000 Allerton Conference*, 2000.
  - [61] I.Ch. Paschalidis and J.N. Tsitsiklis. Congestion-Dependent Pricing of Network Services. *IEEE/ACM Transactions on Networking*, 8(2):171–184, 2000.
  - [62] P. Reichl, P. Flury, J. Gerke, and B. Stiller. How to overcome the feasibility problem for tariffing internet services: the cumulus pricing scheme. In *Proceedings of IEEE ICC 2001, vol. 7*, pages 2079–2083, 2001.
  - [63] P. Reichl and B. Stiller. Edge pricing in space and time: Theoretical and practical aspects of the cumulus pricing scheme. In *Proceedings of the 17th International Teletraffic Congress*, 2001.
  - [64] P. Reichl and B. Stiller. Nil nove sub sole: Why internet charging schemes look like as they do. In *Proceedings of the 4th Berlin Internet Economic Workshop*, 2001.
  - [65] P. Reichl, B. Stiller, and S. Leinen. Auction Models for Multiprovider Internet Connections. In *Proc. Messung, Modellierung und Bewertung MMB'99. Trier (Germany)*, 1999.
  - [66] J.W. Roberts. Quality of Service Guarantees and Charging in Multiservice Networks. *IEICE Trans. Commun.*, E81(5):824–831, 1998.
  - [67] J.W. Roberts and L. Massoulié. Bandwidth sharing and admission control for elastic traffic. *Telecommunication Systems*, 15(1-2):185–201, 2000.
  - [68] N. Semret. *Market Mechanisms for Network Resource Sharing*. PhD thesis, Columbia University, 1999.
  - [69] N. Semret, R.R.-F. Liao, A.T. Campbell, and A.A. Lazar. Market Pricing of Differentiated Internet Services. In *Proceedings of the 7th International Workshop on Quality of Service*, 1999.
  - [70] N. Semret, R.R.-F. Liao, A.T. Campbell, and A.A. Lazar. Pricing, provisioning and peering: Dynamic markets for differentiated internet services

- 
- and implications for network interconnections. *IEEE Journal on Selected Areas in Communications*, 18(12):2499–2513, 2000.
- [71] S. Shenker. Service models and pricing policies for an integrated services internet. In *Performance of "Public Access to the Internet"*, 1993.
- [72] S. Shenker, D. Clark, D. Estrin, and S. Herzog. Pricing in computer networks: reshaping the research agenda. *Computer Communication Review*, 26(2):19–43, 1996.
- [73] D. Songhurst (ed.). *Charging Communication Networks: from Theory to practice*. Elsevier, Amsterdam, 1999.
- [74] B. Stiller, P. Reichl, and S. Leinen. Pricing and Cost Recovery for Internet Services: Practical Review, Classification, and Application of Relevant Models. *Netnomics*, 2(1), 2000.
- [75] B. Tuffin. Revisited Progressive Second Price Auction for Charging Telecommunication Networks. Technical Report 4176, INRIA, May 2001.
- [76] Q. Wang, J.M. Peha, and M.A. Sirbu. Optimal Pricing for Integrated Services Networks. In Lee W. McKnight and Joseph P. Bailey, editors, *Internet Economics*, pages 353–376. MIT Press, 1997.
- [77] H. Yaïche, R.R. Mazumdar, and C. Rosenberg. A Game Theoretic Framwork for Bandwidth Allocation and Pricing in Broadband Networks. *IEEE/ACM Transactions on Networking*, 8(5):667–678, 2000.
- [78] L.S. Zhang, D. Deering, S. Shenker, and D. Zappala. RSVP: A resource ReSerVation Protocol. *IEEE Network Magazine*, 1993.



---

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY  
Unité de recherche INRIA Rennes, Irista, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN  
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

---

Éditeur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399