

Proportional QoS in Differentiated Services Networks: Capacity Management, Equilibrium Analysis and Elastic Demands

Ishai Menache and Nahum Shimkin

Department of Electrical Engineering
Technion, Israel Institute of Technology
Haifa 32000, Israel
{imenache@tx,shimkin@ee}.technion.ac.il

Abstract. The Differentiated Services (Diffserv) architecture is a scalable solution for providing Quality of Service (QoS) over packet switched networks. By its very definition, Diffserv is not intended to provide strict performance guarantees to its subscribers. We propose in this paper a particular form of *relative* performance guarantees. Specifically, the network manager's goal is to maintain pre-defined ratios between common congestion measures over the different service classes. We assume that each service class is advertised with a constant price. Thus, in order to induce its goal, the manager dynamically allocates available capacity between the service classes. This scheme is studied within a network flow model, with self-optimizing users, where each user can choose the amount of flow to ship on each service class according to its service utility and QoS requirements. We pose the entire problem as a *non-cooperative game*. Concentrating on a simplified single-link model with multiple service classes, we establish the existence and uniqueness of the Nash equilibrium where the relative performance goal is obtained. Accordingly, we show how to compute and sustain the required capacity assignment. The extension to a general network topology is briefly outlined.

1 Introduction

Background and Motivation. The need for providing service differentiation over the Internet has been an ongoing concern in the networking community. The Differentiated Services (Diffserv) architecture [5] has been proposed by the IETF as a scalable solution for QoS provisioning. Instead of reserving resources per session (e.g., as in the Integrated Services (IntServ) model [24]), packets are marked to create a smaller number of packet classes, which offer different service qualities. The Diffserv proposal suggests to combine simple priority mechanisms at the network core with admission control mechanisms at the network edges only, in order to create diverse end-to-end services.

The two principal Diffserv classes that have been formalized are the Expedited Forwarding (EF) [8] and the Assured Forwarding (AF) [11] services. The premise of the EF is to provide no-loss and delay reduction to its subscribers.

AF is intended for users who need reliable forwarding even in times of network congestion. Ongoing IETF work concentrates on defining the engineering and architectural aspects of Diffserv-enabled routers (e.g., [2]). However, current technical specifications deliberately do not quantify the actual *service characteristics*, which users will obtain by using the above mentioned classes. Apparently, service characteristics would have to be defined and publicly declared in order to make the distinction between the service classes meaningful to the user (and possibly worth paying for).

The Diffserv network cannot offer strict quality guarantees, as resources are allocated to the service classes based on some average network conditions [7]. Instead, the provider may declare upper-bounds on QoS measures, or alternatively provide looser guarantees, such as probabilistic or time-dependent guarantees. Another option is to offer *relative* quality guarantees. This option can be easily quantified and advertised, as illustrated by the following exemplifying rule: “service class Q will offer an end-to-end average delay, which is at least two times less than any other class, independent of the level of congestion”. When a user buys class Q it is aware of what it gets, and expects the provider to uphold the agreement conditions. In this paper, we focus on the *proportional* QoS model, whereby QoS ratios between the classes are announced. Specifically, we concentrate on *delay* ratios, although analogous definitions may be suggested when considering other QoS measures. The proportional QoS model benefits from implementation-related pros. Ratios are easier to maintain in comparison with absolute end-to-end guarantees, primarily because they may hold for different levels of congestion, and secondly because keeping the ratios locally (on a node basis) leads to fulfilling this objective on the network level.

We shall examine capacity allocation as the main network management tool for achieving the proportional QoS design objective. The goal of capacity allocation is to keep the announced delay ratios, irrespectively of complementary means for network traffic control, such as pricing and admission control. Our focus in this paper is on a simplified single link network, where the network manager owns a fixed amount of capacity to be divided among the link’s offered service classes in order to maintain the QoS ratios objective. Generally, it is easy to calculate the appropriate capacity allocation when the traffic in each service class is fixed. In this paper, however, we consider the interaction of the user behavior and network conditions. Within a standard flow model (see [4], Sec. 5.4) we represent the user population as a finite set of self-optimizing decision makers. The users are heterogenous with respect to their cost functions, which reflect their *price-quantity-quality* tradeoffs. Furthermore, users may modify their flow quantities (i.e., *elastic* users [12]), and may also shift their traffic from one service class to the other in response to current congestion conditions. We pose the overall problem as a non-cooperative game between the manager and the (selfish) users, and explore the associated capacity management policies and equilibrium conditions.

In a recent paper [15], we have considered a similar problem with static (fixed) user demand in the general network context. The present paper extends

this work (for a single-link network) to the case of elastic demand, as well as adding a price term to the overall user's cost.

Related Literature. *Service differentiation approaches:* Several research papers have addressed the model of finitely many classes, in which no strict performance guarantees are given. The simplest approach for providing differentiated services is Odlyzko's Paris Metro Pricing (PMP) proposal [19]. The idea of PMP is to create service differentiation by giving a different price to each service class. Other papers [10, 14, 17] explicitly consider elements such as the user model, the scheduling mechanism and the network objective (e.g., a social or an economic objective). The major concern is usually in calculating the prices that would lead to the network objective. An additional common ground is that the network does not declare any kind of QoS guarantees. Users are assumed to acquire the best deal there is with respect to their quality-price tradeoff. We deviate from the last assumption, by considering the upholding of the service characteristic as a primary management priority.

Selfish routing: Since our model assumes that users are allowed to split traffic between the service classes, our work is related to selfish routing models. Game-theoretic analysis is widely used to study the working conditions of these models. The involved issues include the existence and uniqueness of an equilibrium point, its calculation, and its properties (such as the degradation of performance due to selfish behavior, known as the "price of anarchy" [23]). A common routing model, originated from the field of transportation networks, has considered networks shared by infinitesimal users (see [1] for a survey). The case of finitely many users, each carrying substantial flow has been introduced to the networking literature more recently (see [13, 20, 21]). In [13], the equilibrium properties are applied for network design (namely, link capacity assignment), where the objective is to obtain the socially optimal working point. We use a similar routing model to represent the user's choice of service class.

Proportional QoS: Dovrolis et al. [9] proposed a class of schedulers, based on the Proportional Delay Differentiation (PDD) model. This model aims at providing predetermined delay ratios. The schedulers are implemented by observing the history of the encountered delays (or alternatively, by measuring the delay of the packet at the head of each service class), and serving the class which most deviates from its nominal delay ratio. In the present work we do not rely on PDD schedulers, but rather use a capacitated links model that may be considered a proxy to existing scheduling schemes such as GPS.

Contribution and Organization. This paper proposes schemes for inducing proportional QoS through capacity allocation, when users can react to the allocation decisions. The precise definitions of the network and user models are given in Section 2. The analysis of this model is presented in Section 3, which establishes the existence and uniqueness of the equilibrium point for the network-users game, and presents an algorithm for its computation. An explicit formula is obtained for the best response map of the network, namely the capacity assignment that ensures the QoS-ratio objective for *fixed* network flows, which may be used as a basis for an adaptive capacity assignment scheme. Due to length

constraints, the proofs of some claims are omitted, the reader is referred to [16] for full details.

2 The Single Link Model

Our basic model considers a single link, which supports several service classes. As a stand-alone model, it may be viewed as an approximation of a single path in a network, where the variations in traffic due to other (intersecting) network paths are neglected. Let $\mathcal{I} = \{1, 2, \dots, I\}$ be a finite set of users, which share a link that offers a set of service classes $\mathcal{A} = \{1, 2, \dots, A\}$. Since each service class is characterized by its own price and performance measure (to be described in the sequel), it would be convenient to consider the link with its respective service classes as a two terminal (source-destination) network, which is connected by a set of parallel arcs. Each arc represents a different service class. Thus, the set of arcs is also denoted by \mathcal{A} , and the terms service class and arc are used interchangeably. We denote by f_a^i the flow which user i ships on arc a . User i is free to choose any assignment of $f_a^i \geq 0$. The total demand of each user will be denoted by $f^i \triangleq \sum_{a \in \mathcal{A}} f_a^i$. Turning our attention to an arc $a \in \mathcal{A}$, let f_a be the total flow on that arc, i.e., $f_a = \sum_{i \in \mathcal{I}} f_a^i$. Also, denote by \mathbf{f}_a the vector of all user flows on arc a , i.e., $\mathbf{f}_a = (f_a^1, \dots, f_a^I)$. The user flow configuration \mathbf{f}^i is the vector $\mathbf{f}^i = (f_1^i, \dots, f_A^i)$. The flow configuration \mathbf{f} is the vector of all user flow configurations, $\mathbf{f} = (\mathbf{f}^1, \dots, \mathbf{f}^I)$. A user flow configuration is *feasible* if its components obey the nonnegativity constraints, as described above. We denote by \mathbf{F}^i the set of all feasible user flow configurations \mathbf{f}^i , and by \mathbf{F} the set of all feasible flow configurations \mathbf{f} .

The network manager has a constant capacity C , to be divided between the service classes. This capacity cannot be statically assigned, since the manager cannot predict in advance the number of customers and their preferences. Practically, the network manager would modify the current capacity assignment at slower time scales than the user routing decisions, after periodically measuring class performance. We denote by c_a the allocated capacity at arc a . The capacity allocation of the manager is the vector $\mathbf{c} = (c_1, \dots, c_A)$. An allocation \mathbf{c} is feasible if its components obey the nonnegativity and total capacity constraint, namely (i) $c_a \geq 0$, $a \in \mathcal{A}$ and (ii) $\sum_{a \in \mathcal{A}} c_a = C$. The set of all feasible capacity allocations \mathbf{c} is denoted by Γ . Finally, a system configuration is feasible if it is composed of a feasible flow configuration and a feasible capacity allocation.

Pricing. Each service class $a \in \mathcal{A}$ has a *constant* price p_a per unit traffic. Thus, the network usage price of user i is $\sum_{a \in \mathcal{A}} f_a^i p_a$. Prices may be viewed as an indirect mean for admission control [7], and (among other things) prevent flooding of the better service classes. In this paper, however, we concentrate on capacity assignment as the management tool, assuming that prices are static (or change on a slower time scale). The issue of price setting in our context is left for future work.

The *performance measures* of both the users and the manager are specified through their respective cost functions, which they wish to minimize. We denote

by $J^i(\mathbf{f}, \mathbf{c})$, $i \in \mathcal{I}$ the cost function of user i , and by $J^M(\mathbf{f}, \mathbf{c})$ the cost function of the manager. The costs of the users and the manager are related to the *congestion* level at each of the service classes. We shall use the well known *M/M/1 latency function*

$$D_a(f_a, c_a) = \begin{cases} \frac{1}{c_a - f_a} & f_a < c_a \\ \infty & \text{otherwise} \end{cases} \quad (1)$$

as a congestion measure of each service class. Here c_a is the transmission capacity measured in the same units as the flow f_a . This latency function is often used to model delay costs in communication networks [4], and provides a clear sense of link capacity.

User cost functions. Users are distinguished by their utility function $U^i(f^i)$, which quantifies their utility for shipping a total flow f^i . We make the following assumption on U^i .

Assumption 1 *For every user $i \in \mathcal{I}$, the utility function $U^i : [0, C] \rightarrow \mathfrak{R}$ is increasing, bounded, concave and continuously differentiable.*

We note that utility functions with the above characteristics are commonly used within the networking pricing literature [7, 12]. We define U^i in the set $[0, C]$ only, since a total flow of $f^i > C$ cannot be accommodated by the network. We note that a user may split its flow among different service classes in order to minimize the total cost. The total cost J^i of each user i is comprised from its delay cost, its network usage price, minus its utility, namely

$$J^i(\mathbf{f}, \mathbf{c}) = \beta^i \sum_{a=1}^A f_a^i D_a(f_a, c_a) + \sum_{a=1}^A f_a^i p_a - U^i(f^i), \quad (2)$$

where $\beta^i > 0$ represents the delay sensitivity of user i , and D_a is defined in (1).

Manager cost function. The objective of the manager is to impose certain predetermined ratios between the average delays of the service classes. Taking the delay of class 1 as a reference, the ratios are described by a vector $\rho = (\rho_2, \dots, \rho_A)$, $0 < \rho_a < \infty$, where the manager's objective is to have the delays D_1, \dots, D_A obey

$$D_a(f_a, c_a) = \rho_a D_1(f_1, c_1), \quad (3)$$

where $\rho_1 \triangleq 1$. We refer to that relation as the *fixed ratio objective*. For concreteness, we may assume that $\rho_1 \leq \rho_2 \leq \dots \leq \rho_A$, so that service classes are ordered from best to worst. Similarly, prices are expected to satisfy $p_1 \geq p_2 \geq \dots \geq p_A$, although this is not essential for our results. In functional terms, the manager's cost function may thus be defined as

$$J^M(\mathbf{f}, \mathbf{c}) = \begin{cases} 0 & \text{if (3) holds,} \\ \infty & \text{otherwise.} \end{cases} \quad (4)$$

An alternative objective of the manager that will be considered, is to minimize a weighted sum of the delay functions, that is

$$\bar{J}^M(\mathbf{f}, \mathbf{c}) = \sum_{a \in \mathcal{A}} w_a D_a(f_a, c_a), \quad (5)$$

where $w_a > 0$, $a \in \mathcal{A}$. As we shall see, there is a close relation between the fixed ratios objective and the weighted sum objective (5).

The interaction between the manager and the users will be referred to as the *users-manager game*. A Nash Equilibrium Point (NEP) of that game is a feasible system configuration $(\tilde{\mathbf{f}}, \tilde{\mathbf{c}})$ such that

$$\begin{aligned} J^M(\tilde{\mathbf{f}}, \tilde{\mathbf{c}}) &= \min_{\mathbf{c} \in \Gamma} J^M(\tilde{\mathbf{f}}, \mathbf{c}), \\ J^i(\tilde{\mathbf{f}}^i, \tilde{\mathbf{f}}^{-i}, \tilde{\mathbf{c}}) &= \min_{\mathbf{f}^i \in \mathbf{F}^i} J^i(\mathbf{f}^i, \tilde{\mathbf{f}}^{-i}, \tilde{\mathbf{c}}) \quad \forall i \in \mathcal{I}, \end{aligned} \quad (6)$$

where $\tilde{\mathbf{f}}^{-i}$ stands for the flow configurations of all users, but the i th user. Namely, the NEP is a network working point, where no user, nor the manager, finds it beneficial to change its flow or capacity allocation. Our users-manager game is characterized by the finiteness of the NEP costs, since users can always ship a flow of zero to encounter a finite cost. We shall formally prove this attribute in the next section.

3 Capacity Assignment and Equilibrium Analysis

In this section we analyze the equilibrium point and suggest capacity assignment schemes that induce the ratios objective. First, we show that the manager has a unique best response with respect to the ratio objective (3). This response is a simple solution to a set of linear equations. Then we prove the existence and uniqueness of an equilibrium point, in which the desired ratios are met. Accordingly, we show the equivalence between the best response with respect to (3) and the best response with respect to the modified manager objective function (5) and discuss its implications.

Theorem 1 considers the best response capacity assignment of the manager, given any (fixed) flow configuration.

Theorem 1. *Consider a fixed flow configuration (f_1, \dots, f_A) and a desired ratio vector ρ . If $\sum_{a \in \mathcal{A}} f_a < C$, there exists a unique capacity allocation $\mathbf{c} \in \Gamma$ such that (3) is met. This allocation is explicitly given by*

$$c_a - f_a = (C - \sum_{\alpha \in \mathcal{A}} f_\alpha) \frac{\rho_a^{-1}}{\sum_{\alpha \in \mathcal{A}} \rho_\alpha^{-1}}. \quad (7)$$

Proof. The allocation (7) is derived from solving the following set of *linear* equations: $\rho_a(c_a - f_a) = (c_1 - f_1)$, $a = 2, \dots, A$; and $\sum_{a=1}^A c_a = C$. \square

The above result allows the manager to explicitly calculate its *best response* assignment, namely the capacity assignment that will satisfy the fixed ratio objective given the *current* network flows. This calculation requires just the total flows in each service class, which are easy to measure. The next theorem establishes the existence and uniqueness of the equilibrium point.

Theorem 2. *For every delay ratios vector ρ , there exists a unique Nash equilibrium point. This NEP has finite costs for the manager and for the users. In particular, the ratios objective of the manager is satisfied.*

Proof. The proof follows from the next four lemmas.

Lemma 1. *For every delay ratio vector ρ , there exists a NEP. Furthermore, in every NEP the costs of the users and the manager are finite and the ratio objective is met.*

Proof. See Appendix A.2.

Lemma 2. *Let D_1, \dots, D_A be the class delays at some NEP. Then the following equations are met at the equilibrium for every $i \in \mathcal{I}$ and every $a \in \mathcal{A}$*

$$\begin{aligned} \beta^i (D_a + f_a^i D_a^2) + p_a &= U^{i'}(f^i) \quad \text{if } f_a^i > 0, \\ \beta^i D_a + p_a &\geq U^{i'}(f^i) \quad \text{if } f_a^i = 0. \end{aligned} \quad (8)$$

Proof. Observe that

$$\begin{aligned} \frac{dJ^i(\mathbf{f}, \mathbf{c})}{df_a^i} &= \beta^i \left(\frac{1}{(c_a - f_a)} + \frac{f_a^i}{(c_a - f_a)^2} \right) + p_a - U^{i'}(f^i) \\ &= \beta^i (D_a(f_a, c_a) + f_a^i D_a^2(f_a, c_a)) + p_a - U^{i'}(f^i). \end{aligned} \quad (9)$$

Then (8) may be readily seen to be the KKT optimality conditions [6] for minimizing the cost function (2) of user i subject to the flow constraint $f_a^i \geq 0$. \square

Lemma 3. *Consider a NEP with given class delays D_1, \dots, D_A . Then the respective equilibrium flows f_a^i are uniquely determined.*

Proof. Consider the following optimization problem in (f_1^i, \dots, f_A^i) :

$$\begin{cases} \min \sum_{a=1}^A \frac{1}{2} \beta^i D_a^2 f_a^i{}^2 + f_a^i (\beta^i D_a + p_a) - U^i(\sum_a f_a^i) \\ \text{s.t.} \quad f_a^i \geq 0 \end{cases}, \quad (10)$$

where we assume that the delays D_a are fixed. Note that (10) is a strictly convex optimization problem, since the objective function is the sum of a diagonal quadratic term (with $\beta^i D_a^2 > 0$ for every a) and the negation of U^i , where U^i is concave by Assumption 1. Thus, this problem has a unique minimum, which is characterized by the KKT optimality conditions. It is now readily seen that the KKT conditions for (10) coincide with the conditions in (8). Thus, by Lemma 2, any set of equilibrium flows $(f_a^i)_{a \in \mathcal{A}}$ is a solution of (10). But since this solution is unique, the claim is established. \square

Lemma 4. *Consider two Nash equilibrium points (\mathbf{f}, \mathbf{c}) and $(\tilde{\mathbf{f}}, \tilde{\mathbf{c}})$. Then $D_a(f_a) = D_a(\tilde{f}_a)$ for every $a \in \mathcal{A}$.*

Proof. Define $D_a \triangleq D_a(f_a)$ and $\tilde{D}_a \triangleq D_a(\tilde{f}_a)$. Assume that

$$\tilde{D}_\alpha > D_\alpha \text{ for some } \alpha \in \mathcal{A}. \quad (11)$$

Then $\tilde{D}_a > D_a$ for every $a \in \mathcal{A}$ since the ratios are met in both equilibria (Lemma 1). Noting (4) and (7) it follows from (11) that $\sum_{a \in \mathcal{A}} \tilde{f}_a^j > \sum_{a \in \mathcal{A}} f_a^j$, which implies that there exists some user j for which

$$\tilde{f}^j = \sum_{a \in \mathcal{A}} \tilde{f}_a^j > \sum_{a \in \mathcal{A}} f_a^j = f^j. \quad (12)$$

We next contradict (12) by invoking the next two implications:

$$f_a^j = 0 \Rightarrow \tilde{f}_a^j = 0 \quad (13)$$

$$f_a^j > 0 \Rightarrow \tilde{f}_a^j > f_a^j. \quad (14)$$

Their proof is based on the KKT conditions (8). Since the utility U^j is concave, then by (12) we have $\lambda^j \triangleq U^{j'}(f^j) \geq U^{j'}(\tilde{f}^j) \triangleq \tilde{\lambda}^j$. If $f_a^j = 0$, then $\beta^j D_a + p_a \geq \lambda^j \geq \tilde{\lambda}^j$. Since $\tilde{D}_a > D_a$, then $\beta^j \tilde{D}_a + p_a > \lambda^j$, hence $f_a^j = 0$. To prove (14) note first that it holds trivially if $\tilde{f}_a^j = 0$. Next assume $\tilde{f}_a^j > 0$. Then by (8)

$$\beta^j D_a + \beta^j D_a^2 f_a^j + p_a = \lambda^j \geq \tilde{\lambda}^j = \beta^j \tilde{D}_a + \beta^j \tilde{D}_a^2 \tilde{f}_a^j + p_a. \quad (15)$$

Since $\tilde{D}_a > D_a$ (hence $\tilde{D}_a^2 > D_a^2$), and β^j are positive, we must have $f_a^j > \tilde{f}_a^j$ in order for (15) to hold, which establishes (14). Summing user j 's flows according to (13)-(14) yields $\sum_{a \in \mathcal{A}} \tilde{f}_a^j \leq \sum_{a \in \mathcal{A}} f_a^j$, which contradicts (12). Thus $\tilde{D}_a \leq D_a$. Symmetrical arguments will lead to $\tilde{D}_a \geq D_a$, i.e., $\tilde{D}_a = D_a$, hence $\tilde{D}_a = D_a$ for every $a \in \mathcal{A}$. \square

The last two lemmas imply that the user flows and the class delays in equilibrium are unique. The capacities in the equilibrium must also be unique since $c_a = f_a + \frac{1}{D_a}$, where $f_a = \sum_{i \in I} f_a^i$. This establishes the uniqueness of the NEP, and completes the proof of Theorem 2. \square

A possible criticism of the fixed ratio objective, as defined in (3), is that it does not account at all for absolute congestion measures, namely the delays themselves rather than their ratios. However, the next result shows that by achieving the ratio objective, the manager in fact minimizes the cost function \bar{J}^M (defined in (5)), which is just an appropriately weighted sum of the delays over the different service classes.

Theorem 3. *Consider the users-manager game, whose NEP is defined in (6), and an additional game, which is similar except that J^M is replaced by \bar{J}^M . If the parameters are such that $w_a = \frac{1}{\rho_a^2}$ for every $a \in \mathcal{A}$, where $\rho_1 \triangleq 1$, then the two games are equivalent in the sense that their (unique) equilibrium points coincide.*

Proof. See Appendix A.1.

Note that the weights w_a are inversely proportional to ρ_a^2 , which assigns higher weight to better (lower relative delay) service classes, as might be expected. From an algorithmic point of view, it should be noted that \bar{J}^M is a convex and continuous function, and therefore may provide a sound basis for an iterative (e.g., gradient-based) algorithm that may be used by the manager to

minimize this cost, with the goal of eventually reaching the desired fixed-ratio equilibrium. Yet, a specific consideration of such an algorithm is beyond the scope of this work.

We conclude this section by considering the computation of the NEP. Recall that the uniqueness of the user flows was established in Lemma 3 via the definition of strictly convex optimization problems. Hence, by solving the same optimization problems, the NEP can be efficiently calculated. The only unknown variable which is required for the computation is D_1 . In our case, an iterative search for D_1 may be easily performed, by comparing the total flow $\sum f_a$ (for a given D_1), which is obtained from both the best response formula (7) and also from the solutions to (10). More details on the search method are given in Appendix B.

Remark 1. General Networks. It is obviously of interest to extend the results of this section to a general network topology. In a recent paper [15], we have made such an extension for a model with fixed (plastic) user demands. Since the same extension can be similarly used here, we briefly outline its key features. The following setup is applied: (i) each user has a unique fixed path from its source to its destination; (ii) the QoS ratio objective is maintained on a link basis, i.e., management is performed via a *distributed* approach, where the capacity adaptation is performed locally, at the link level. Observe that if the above two features are maintained, then the QoS ratios are met *end-to-end* for every user. The game framework now includes network *managers*, one for every link. Due to the locality of the capacity management, the formula for the best response map (Theorem 1) and the use of \bar{J}^M instead of J^M (Theorem 3) are trivially extended to the general network case; so is the proof for the existence of the equilibrium. The issue of uniqueness of the equilibrium (as in Lemma 4) is however more complicated in the general network case and is currently an open problem.

4 Conclusion

The proposed approach to QoS provisioning in Diffserv networks focuses on maintaining relative congestion measures in the different service classes. Our analysis demonstrates the feasibility of this approach, and in particular establishes the existence and uniqueness of a working point, which satisfies this QoS objective in an elastic, reactive and heterogenous user environment. Our results provide an efficient algorithm for computing the Nash equilibrium. However, from the manager's viewpoint, this computation requires complete knowledge of the users' preferences, which may not be available. An alternative scheme is to use adaptive capacity assignment, for example by utilizing the best-response map (7), which requires only the total flows in each service class, which are easy to measure. The analysis of such a scheme is an important issue for future research. Additional research topics include the price setting issue, proportional QoS with other cost functions and QoS measures, and the equilibrium dynamics of capacity allocation algorithms.

References

1. E. Altman and L. Wynter. Equilibrium, games, and pricing in transportation and telecommunications networks. *Networks and Spatial Economics*, 4(1):7–21, 2004.
2. F. Baker, K. Chan, and A. Smith. Management information base for the differentiated services architecture. RFC 3289, 2002.
3. T. Basar and G. J. Olsder. *Dynamic Noncooperative Game Theory*. Academic Press, 1995.
4. D. P. Bertsekas and R. G. Gallager. *Data Networks*. Prentice-Hall, 1992.
5. S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss. An architecture for differentiated services. RFC 2475, 1998.
6. S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2003.
7. C. Courcoubetis and R. Weber. *Pricing Communication Networks: Economics, Technology and Modelling*. Wiley, 2003.
8. B. Davie, A. Charny, J. Bennett, K. Benson, J. L. Boudec, W. Courtney, S. Davari, V. Firoiu, and D. Stiliadis. An expedited forwarding PHB (per-hop behavior). RFC 3246, 2001.
9. C. Dovrolis, D. Stiliadis, and P. Ramanathan. Proportional differentiated services: Delay differentiation and packet scheduling. *IEEE/ACM Transactions on Networking*, 2002.
10. Y. Hayel, D. Ros, and B. Tuffin. Less-than-best-effort services: Pricing and scheduling. In *Proceedings of IEEE INFOCOM*, 2004.
11. J. Heinanen, F. Baker, W. Weiss, and J. Wroclawski. Assured forwarding PHB group. RFC 2597, 1999.
12. 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*, 49:237–252, 1998.
13. Y. A. Korilis, A. A. Lazar, and A. Orda. Architecting noncooperative networks. *IEEE Journal on Selected Areas in Communication*, 13(7):1241–1251, 1995.
14. M. Mandjes. Pricing strategies under heterogeneous service requirements. In *Proceedings of IEEE INFOCOM*, pages 1210–1220, 2003.
15. I. Menache and N. Shimkin. Capacity assignment in Diffserv networks for proportional QoS. Technical Report CCIT-544, Technion, July 2005.
16. I. Menache and N. Shimkin. Proportional QoS in differentiated services networks: Capacity management, equilibrium analysis and elastic demands. Available from <http://www.ee.technion.ac.il/people/shimkin/preprints/MS05x.pdf>, July 2005. Extended version.
17. H. Mendelson and S. Whang. Optimal incentive-compatible priority pricing for the M/M/1 queue. *Operations Research*, 38:870–883, 1990.
18. J. Nocedal and S. J. Wright. *Numerical optimization*. Springer, 1999.
19. A. M. Odlyzko. Paris metro pricing for the internet. In *Proceedings of the ACM Conference on Electronic Commerce*, pages 140–147, 1999.
20. A. Orda, R. Rom, and N. Shimkin. Competitive routing in multi-user environments. *IEEE/ACM Transactions on Networking*, 1:510–521, 1993.
21. O. Richman and N. Shimkin. Topological uniqueness of the Nash equilibrium for non-cooperative routing. *Mathematics of Operations Research*, submitted for publication, October 2004.
22. J. B. Rosen. Existence and uniqueness of equilibrium points for concave n-person games. *Econometrica*, 33(3):520–534, 1965.

23. T. Roughgarden. *Selfish Routing and the Price of Anarchy*. MIT Press, 2005.
24. X. Xiao and L. M. Ni. Internet QoS: A big picture. *IEEE Network*, 36(11):8–18, 1999.

APPENDIX

A Proofs

A.1 Proof of Theorem 3

The proof is based on characterizing the manager's best response. The characterization is established in the following lemma.

Lemma 5. *Consider a fixed flow vector (f_1, \dots, f_A) , $\sum_{a \in \mathcal{A}} f_a < C$, and a manager whose cost function is given by (5). Then $\mathbf{c} = (c_1, \dots, c_A)$ is the manager's best response to (5) if and only if the following relations hold for every $a, a' \in \mathcal{A}$:*

$$\sqrt{w_a} D_a(c_a, f_a) = \sqrt{w_{a'}} D_{a'}(c_{a'}, f_{a'}). \quad (16)$$

Proof. Observe first that the manager's optimal capacity allocation must obey $0 < c_a < C$ for every $a = 1, \dots, A$. Indeed, if $c_a = 0$ for some $a \in \mathcal{A}$ then the manager cost is infinite, as the summand $w_a D_a(f_a, 0)$ in (5) is infinite; if $c_a = C$ for some $a \in \mathcal{A}$ then the manager cost is also infinite, as all summands $w_\alpha D_\alpha(f_\alpha, c_\alpha)$ $\alpha \in \mathcal{A}, \alpha \neq a$ in (5) are infinite. Thus, the only active constraint of the manager is $\sum_{a \in \mathcal{A}} c_a = C$. Consequently, the KKT optimality conditions for minimizing the cost function \bar{J}^M (which are necessary and sufficient in our case, since $w_a D_a(f_a, c_a)$ is convex in c_a) imply that $\frac{d\bar{J}^M}{dc_a} = \frac{d\bar{J}^M}{dc_{a'}}$ for every $a, a' \in \mathcal{A}$ (this equality is due to the equality of each $\frac{d\bar{J}^M}{dc_a}$ to the Lagrange multiplier of the active constraint). Noting that $\frac{d\bar{J}^M}{dc_a} = \frac{-w_a}{(c_a - f_a)^2}$, we get

$$\frac{w_a}{(c_a - f_a)^2} = \frac{w_{a'}}{(c_{a'} - f_{a'})^2}. \quad (17)$$

Taking a square root from both sides of the last equation (and noting that $c_a - f_a > 0$ for every $a \in \mathcal{A}$) gives the required relations (16). \square

Now substituting $w_a = \frac{1}{\rho_a^2}$ for every $a = 1, \dots, A$, $\rho_1 \triangleq 1$ in (16) yields the delays whose ratios are described by ρ . These ratios are obviously the manager's best response to (f_1, \dots, f_A) when using J^M with ρ as the desired ratio vector. Thus, the best responses of J^M with ρ and \bar{J}^M with $\mathbf{w} = (w_1, \dots, w_A)$ coincide. Since the user's cost function is given by (2) in both users-manager games, the Nash equilibrium points of both games coincide. \square

A.2 Proof of Lemma 1

As in Rosen [22], we shall apply the Kakutani fixed point theorem for the proof. Hence we first state it precisely, along with the necessary mathematical definitions. These are taken from [3].

Definition 1. (*Upper semicontinuity*) Let Λ be a function defined on a normed linear space X , and associating with each $x \in X$ a subset $\Lambda(x)$ of some (other) normed linear space Y . Then, Λ is said to be upper semicontinuous (usc) at a point $x_0 \in X$, if for any sequence $\{x_i\}$ converging to x_0 and any sequence $\{y_i \in \Lambda(x_i)\}$ converging to y_0 , we have $y_0 \in \Lambda(x_0)$. The function Λ is upper semicontinuous if it is usc at each point of X .

Theorem 4. (*Kakutani*) Let S be a compact and convex subset of \mathbb{R}^n . and let Λ be an upper semicontinuous function which assigns to each $x \in S$ a closed and convex subset of S . Then there exists some $x \in S$ such that $x \in \Lambda(x)$.

The proof of Lemma 1 is obtained by the following steps:

1. *Definition of S .* Let r^1, \dots, r^I be (arbitrary) positive constants such that $r^i > C$ for every $i \in \mathcal{I}$. Define $S \subset \mathbf{F} \times \mathbf{\Gamma}$ to be the following compact and convex product set

$$S = \{(\mathbf{f}, \mathbf{c}) \in S : \sum_{a \in \mathcal{A}} f_a^i \leq r^i, f_a^i \geq 0 \forall i; \sum_{a \in \mathcal{A}} c_a = C, c_a \geq 0\}. \quad (18)$$

Note that a NEP cannot exist outside S , since in case that the flow exceeds the total capacity (which is the case for $(\mathbf{f}, \mathbf{c}) \notin S$), there exists a user with infinite cost (due to the infinite delay cost) which can improve its cost (make it finite) by shipping, e.g., a zero flow to all service classes.

2. *Definition of Λ .* We define the point-to-set mapping $(\mathbf{f}, \mathbf{c}) \in S \rightarrow \Lambda(\mathbf{f}, \mathbf{c})$, as follows.

$$\Lambda(\mathbf{f}, \mathbf{c}) = \{(\hat{\mathbf{f}}, \hat{\mathbf{c}}) \in S : \hat{\mathbf{f}}^i \in \underset{\tilde{\mathbf{f}}^i \in \mathbf{F}^i}{\operatorname{argmin}} J^i(\tilde{\mathbf{f}}^i, \mathbf{f}^{-i}, \mathbf{c}) \forall i \in \mathcal{I}, \hat{\mathbf{c}} = \underset{\tilde{\mathbf{c}} \in \mathbf{\Gamma}}{\operatorname{argmin}} \bar{J}^M(\mathbf{f}, \tilde{\mathbf{c}})\}. \quad (19)$$

Note that for every $(\mathbf{f}, \mathbf{c}) \in S$ we have $(\mathbf{f}, \mathbf{c}) \rightarrow \Lambda(\mathbf{f}, \mathbf{c}) \subset S$ from the same reasoning that was used in the definition of S .

3. *Upper semicontinuity of Λ .*
 - It is straightforward that Λ is usc for the points $(\mathbf{f}, \mathbf{c}) \in S$ such that $\sum_a f_a < C$. Indeed, when $\sum_a f_a < C$, the best response map of $J^i(\mathbf{f}, \mathbf{c})$ and $\bar{J}^M(\mathbf{f}, \mathbf{c})$ is continuous (by the continuity of the cost functions), and finite. In fact, by the strictly convexity property of the cost functions in their respective decision variables, Λ is a point-to-point mapping, in which continuity holds element-wise and thus point-wise.

- For the case where $\sum_a f_a \geq C$, observe that the user’s best response still maintains the continuity and finiteness properties. As to the manager, let $\{(\mathbf{f}^k, \mathbf{c}^k)\}$ be a sequence converging to (\mathbf{f}, \mathbf{c}) such that $\sum_a f_a \geq C$. The manager is indifferent as to its “best response” to \mathbf{f} (since the manager will obtain an infinite cost regardless of the chosen capacity allocation), thus $\{\tilde{\mathbf{c}} | (\tilde{\mathbf{f}}, \tilde{\mathbf{c}}) \in \Lambda(\mathbf{f}, \mathbf{c})\} = \Gamma$, which implies that for any sequence $\{\Lambda(\mathbf{f}^k, \mathbf{c}^k)\}$ converging to some $(\mathbf{f}_0, \mathbf{c}_0)$, we have that $(\mathbf{f}_0, \mathbf{c}_0) \in \Lambda(\mathbf{f}, \mathbf{c})$.
- 4. $\Lambda(\mathbf{f}, \mathbf{c})$ is closed and convex. In case that $\sum_a f_a < C$, we have a point-to-point mapping which is also closed and convex. For the case where $\sum_a f_a \geq C$, $\Lambda(\mathbf{f}, \mathbf{c})$ is a set of points $(\tilde{\mathbf{f}}, \tilde{\mathbf{c}})$, where $\tilde{\mathbf{f}}$ is uniquely determined and $\tilde{\mathbf{c}} \in \Gamma$. This is obviously a close and convex set.
- 5. *Replacing \bar{J}^M with J^M .* We may replace \bar{J}^M by J^M as described in Theorem 3. This replacement is justified by the equivalence of the best response of both functions, under proper weight setting of the latter, as described in Theorem 3.
- 6. *Finiteness of the NEP.* Note that for every system configuration (\mathbf{f}, \mathbf{c}) , if not all costs are finite, then at least one game-player (user or manager) with infinite cost can change its own flow configuration to make its cost finite. This argument is true, since the users can always change their flow allocation to make their cost finite. For the case where the manager is the only game-player with infinite cost, we have $\sum_a f_a < C$ (as user costs are finite). For this case, the manager can apply its best response (7) to obtain a finite cost.
- 7. *Applying the Kakutani fixed point theorem.* Applying the Kakutani fixed point theorem with the above definitions of S and Λ , we conclude that there exists a NEP. This NEP is finite as shown above, and it is also a NEP where the ratios are met (by a proper replacement of \bar{J}^M with J^M).

□

B On the computation of the Nash equilibrium

In this section, we describe how to efficiently calculate the NEP (6). Given D_1 , the equilibrium is calculated via the solution of I optimization problems (10) of A variables each. Since each problem is convex, it may be solved in polynomial time by non-linear optimization techniques [6]. The only issue that needs to be resolved is the calculation of D_1 . We next describe how to derive this (scalar) value, via a standard search procedure. The following observations are needed for the establishment of the search procedure.

1. We conclude from the manager’s best response map (7) that the quantity $\sum_a f_a$ strictly increases with D_1 . This follows from taking the inverse of both sides of (7), followed by simple algebraic operations:

$$\sum_{a \in \mathcal{A}} f_a = C - \frac{\sum_{a \in \mathcal{A}} \rho_a^{-1}}{D_1}. \quad (20)$$

2. The same quantity *strictly* increases with D_1 as the (aggregate) solution to the user optimization problems (10). This fact was formally established in Lemma 4.
3. For the sake of exposition, let us denote the total flow which is obtained from (20) (for a fixed D_1) as f_{MBE} and the total flow which is obtained from (10) (for the same fixed D_1) as f_{UBE} (MBE and UBE stand for manager best response and user best response, respectively). Since f_{MBE} is strictly increasing in D_1 and f_{UBE} is strictly decreasing in D_1 , then there is a unique value of D_1 where these quantities have the same value. The value of D_1 which equalizes f_{MBE} and f_{UBE} is obviously the value of D_1 at the NEP.
4. Following the last argument, we observe that $f_{UBE} - f_{MBE}$ monotonously decreases with D_1 . We may thus use standard search techniques [6, 18], such as the bisection method, in order to obtain the scalar value of D_1 for which $f_{UBE} - f_{MBE} = 0$.