

On Fair and Efficient Bandwidth Allocation by the Multiple Target Approach

Włodzimierz Ogryczak
Warsaw University of Technology
Institute of Control & Comput. Engg.
00-665 Warsaw, Poland
E-mail: ogryczak@ia.pw.edu.pl

Marcin Milewski
Warsaw University of Technology
Institute of Telecommunications
00-665 Warsaw, Poland
E-mail: mmilewsk@elka.pw.edu.pl

Adam Wierzbicki
Polish-Japanese Institute
of Information Technology
02-008 Warsaw, Poland
E-mail: adamw@pjwstk.edu.pl

Abstract—Expanding demand on the Internet services leads to an increased role of the network dimensioning problem for elastic traffic where one needs to allocate bandwidth to maximize service flows and simultaneously to reach a fair treatment of all the elastic services. Thus, both the overall efficiency (throughput) and the fairness (equity) among various services are important. The Max-Min Fairness (MMF) approach, widely used to this problem, guarantees fairness but may lead to significant losses in the overall throughput of the network. In this paper we show how the concepts of multiple criteria equitable optimization can be effectively used to generate various fair resource allocation schemes. We introduce a multiple target model equivalent to equitable optimization and we develop a corresponding procedure to generate fair efficient bandwidth allocations. The procedure is tested on a sample network dimensioning problem and its abilities to model various preferences are demonstrated.

I. INTRODUCTION

Expanding demand on the Internet services has led to an increased role of the traffic carried by the IP protocol in telecommunication networks. Due to the use of packet switching, the IP protocol can provide greater network utilization (the so-called multiplexing gain). For these reasons, network management may be interested in designing networks which allow to extend throughput for the IP protocol. The TCP protocol is the most frequently used transport protocol in best-effort IP networks. The data traffic carried by the TCP protocol adapts its throughput to the amount of available bandwidth. Such traffic, called *elastic traffic*, is capable to use the entire available bandwidth, but it will also be able to reduce its throughput in the presence of contending traffic. Nowadays, the network management often faces the problem of designing networks that carry elastic traffic. These network design problems are, essentially, the network dimensioning problems as they can be reduced to a decision about link capacities. Flow sizes are outcomes of the design problem, since the flows adapt to given network resources on a chosen path.

A straightforward network dimensioning with elastic traffic could be thought of as a search for such network flows that will maximize the aggregate network throughput while staying within a budget constraint for the costs of link bandwidth. However, maximizing aggregate throughput can result in extremely unfair solutions allowing even for starvation of flows

for certain services. On the other extreme, while looking at the problem from the perspective of a network user, the network flows between different nodes should be treated as fairly as possible [2]. Actually, a fair way of distribution of the bandwidth among competing demands becomes a key issue in computer networks [3] and telecommunications network design, in general [16], [18]. The so-called Max-Min Fairness (MMF) [1], [4], [9] is widely considered as such ideal fairness criteria. Indeed, the lexicographic max-min optimization used in the MMF approach generalizes equal sharing at a single link bandwidth to any network while maintaining the Pareto optimality. Certainly, allocating the bandwidth to optimize the worst performances may cause a large worsening of the overall throughput of the network. Therefore, network management must consider two goals: increasing throughput and providing fairness.

The purpose of this work is to show that there exists a multiple criteria model that allows to represent consistently the overall efficiency and fairness goals. Moreover, the criteria measure actual network throughput for various levels (targets) of flows. Thereby, the criteria can easily be introduced into the model and they allow to apply effectively the reference point methodology where the decision maker specifies preferences in terms of aspiration levels (reference point), i.e., by introducing desired (acceptable) levels for several criteria.

The paper is organized as follows. In the next section we formalize the network dimensioning problem we consider. In Section III, basic fair solution concepts for resource allocation are related to the multiple criteria equitable optimization theory and the multiple target model is introduced. In Section IV, the reference point methodology is applied to the multiple target problem allowing us to model various fair and efficient allocation schemes with simple control parameters. Finally, in Section V, we present some results of our initial computational experience with this new approach.

II. THE BANDWIDTH ALLOCATION PROBLEM

The problem of network dimensioning with elastic traffic can be formulated basically as a Linear Programming (LP) based resource allocation model as follows [16]. Given a network topology $G = \langle V, E \rangle$, consider a set of pairs of nodes as the set $I = \{1, 2, \dots, m\}$ of services representing the

elastic flow from source v_i^s to destination v_i^d . For each service, we have given the set P_i of possible routing paths in the network from the source to the destination. This information can be summarized in the form of binary matrices $\Delta_e = (\delta_{eip})_{i \in I, p \in P_i}$ assigned to each link $e \in E$, where $\delta_{eip} = 1$ if link e belongs to the routing path $p \in P_i$ (connecting v_i^s with v_i^d) and $\delta_{eip} = 0$ otherwise.

For each service $i \in I$, the elastic flow from source v_i^s to destination v_i^d is a variable representing the model outcome and it will be denoted by x_i . This flow may be realized along various paths $p \in P_i$. The flow may be either split among several paths or a single path must be finally selected to serve the entire flow. Actually, the latter case of nonbifurcated flows is more commonly required and our analysis is focused on this case. Both bifurcated or nonbifurcated flows may be modeled as $x_i = \sum_{p \in P_i} x_{ip}$ where x_{ip} (for $p \in P_i$) are nonnegative variables representing the elastic flow from source v_i^s to destination v_i^d along the routing path p . Although, the single-path model requires additional multiple choice constraints to enforce nonbifurcated flows.

The network dimensioning problem depends on allocating the bandwidth to several links in order to maximize flows of all the services (demands). Typically, the network is already operated and some bandwidth is already allocated (installed) and decisions are rather related to the network expansion. Therefore, we assume that each link $e \in E$ has already capacity a_e while decision variables ξ_e represent the bandwidth newly allocated to link $e \in E$ thus expanding the link capacity to $a_e + \xi_e$. Certainly, all the decision variables must be nonnegative: $\xi_e \geq 0$ for all $e \in E$ and there are usually some bounds (upper limits) on possible expansion of the links capacities: $\xi_e \leq \bar{a}_e$ for all $e \in E$. Finally, the following constraints must be fulfilled:

$$0 \leq x_{ip} \leq M u_{ip}, \quad u_{ip} \in \{0, 1\} \quad \forall i \in I; p \in P_i \quad (1)$$

$$\sum_{p \in P_i} u_{ip} = 1 \quad \forall i \in I \quad (2)$$

$$\sum_{i \in I} \sum_{p \in P_i} \delta_{eip} x_{ip} \leq a_e + \xi_e \quad \forall e \in E \quad (3)$$

$$0 \leq \xi_e \leq \bar{a}_e \quad \forall e \in E \quad (4)$$

$$\sum_{p \in P_i} x_{ip} = x_i \quad \forall i \in I \quad (5)$$

$$\sum_{e \in E} c_e \xi_e \leq B \quad (6)$$

where (1)–(2) represent single-path flow requirements using additional binary (flow assignment) variables u_{ip} equal 1 if path $p \in P_i$ is assigned to serve flow x_i and 0 otherwise, and a large constant M upper bounding the largest possible total flows x_i . Next, (5) define the total service flows, while (3)–(4) establish the relation between service flows and links bandwidth. The quantity $y_e = \sum_{i \in I} \sum_{p \in P_i} \delta_{eip} x_{ip}$ is the load of link e and it cannot exceed the available link capacity. Further, while allocating the bandwidth to several links in the network dimensioning process the decisions must keep the cost within available budget B for all link bandwidths. This

represented with inequality (6) where for each link $e \in E$ the cost of allocated bandwidth is c_e . In the basic model of network dimensioning it is assumed that any real amount of bandwidth may be installed and marginal costs c_e of link bandwidth is given.

The model constraints (1)–(6) define a Mixed Integer LP (MILP) feasible set. In the simplified problem with linear link dimensioning function ($a_e = 0$ for all links), the cost of the entire path p for service i can be directly expressed by the formula: $\kappa_{ip} = \sum_{e \in E} c_e \delta_{eip}$ and the cheapest path for each service can be then easily identified and preselected. Having preselected routing path for each demand ($|P_i| = 1$) one may consider variable x_i directly as flow along the corresponding path ($x_i = x_{i1}$). Constraints (6) and (3) may be then treated as equations and allowing one to eliminate variables ξ_e , thus formulating the problem as a simplified resource allocation model with only one constraint: $\sum_{i=1}^m \kappa_i x_i = B$ and variables x_i representing directly the decisions. Note that one cannot define directly any cost κ_{ip} of the path $p \in P_i$ when some capacity is already available ($a_e > 0$ for some $e \in E$). In other words, in the problem we consider the cost of available link capacity is actually nonlinear (piecewise linear) and this results in the lack of direct formula for the path cost since it depends on possible sharing with other paths of the preinstalled bandwidth (free capacity a_e).

The network dimensioning model can be considered with various objective functions, depending on the chosen goal. One may consider two extreme approaches. The first extreme is the maximization of the total throughput (the sum of flows) $\sum_{i \in I} x_i$. On the other extreme, the network flows between different nodes should be treated as fairly as possible which leads to the maximization of the smallest flow or rather to the lexicographically expanded max-min optimization (the so-called max-min ordering) allowing also to maximize the second smallest flows provided that the smallest remain optimal, the third smallest, etc. This approach is widely recognized in networking as the so-called Max-Min Fairness (MMF) [1], [4] and it is consistent with the Rawlsian theory of justice [17]. The throughput maximization can always result in extremely unfair solutions allowing even for starvation of certain flows while the MMF solution may cause a large worsening of the throughput of the network. In an example built on the backbone network of a Polish ISP, it turned out that the throughput in a perfectly fair solution could be less than 50% of the maximal throughput [13].

Network management may be interested in seeking a compromise between the two extreme approaches discussed above. One possible approach depends on maximization of the sum of the flows evaluated with some (concave) utility function $\sum_{i \in I} U_i(x_i)$ [11]. However, such an approach requires to build (or to guess) a utility function prior to the analysis and later it gives only one possible compromise solution. It is very difficult to identify the preferences at the beginning of the decision process. Moreover, all the utility functions that really take into account any fairness preferences are nonlinear, thus resulting in computationally hard optimization problems

when applied to the MILP models. In the following, we shall describe an approach that allows to search for such compromise solutions with multiple linear criteria rather than the use nonlinear objective functions. All these criteria represent partial throughput for several target levels of flows.

III. FAIR ALLOCATIONS AND PARTIAL THROUGHPUTS

The bandwidth allocation problem we consider may be viewed as a special case of general resource allocation problem where a set I of m services is considered and for each service $i \in I$, its measure of realization x_i is a function $x_i = f_i(\xi)$ of the allocation pattern $\xi \in A$. This function, called the individual objective function, represents the outcome (effect) of the allocation pattern for service i . In applications, we consider, the measure expresses the service flow and a larger value of the outcome means a better effect (higher service quality or client satisfaction). This leads us to a vector maximization problem:

$$\max \{(x_1, x_2, \dots, x_m) : \mathbf{x} \in Q\} \quad (7)$$

where $Q = \{(x_1, \dots, x_m) : x_i = f_i(\xi), i \in I, \xi \in A\}$ denotes the attainable set for outcome vectors \mathbf{x} . For the network dimensioning problems, we consider, the set Q is an MILP feasible set defined by basic constraints (1)–(6).

Model (7) only specifies that we are interested in maximization of all outcomes x_i for $i \in I$. In order to make it operational, one needs to assume some solution concept specifying what it means to maximize multiple outcomes. The solution concepts are defined by properties of the corresponding preference model within the outcome space. The commonly used concept of the Pareto-optimal solutions, as feasible solutions for which one cannot improve any criterion without worsening another, depends on the rational dominance \succeq_r which may be expressed in terms of the vector inequality: $\mathbf{x}' \succeq_r \mathbf{x}''$ iff $x'_i \geq x''_i$ for all $i \in I$.

In order to ensure fairness in a system, all system entities have to be equally well provided with the system's services. This leads to concepts of fairness expressed by the equitable rational preferences [12], [6]. First of all, the fairness requires impartiality of evaluation, thus focusing on the distribution of outcome values while ignoring their ordering. That means, in the multiple criteria problem (7) we are interested in a set of outcome values without taking into account which outcome is taking a specific value. Hence, we assume that the preference model is impartial (anonymous, symmetric). In terms of the preference relation it may be written as the following axiom

$$(x_{\tau(1)}, x_{\tau(2)}, \dots, x_{\tau(m)}) \cong (x_1, x_2, \dots, x_m) \quad (8)$$

for any permutation τ of I . Further, fairness requires equitability of outcomes which causes that the preference model should satisfy the (Pigou–Dalton) principle of transfers:

$$\mathbf{x} - \varepsilon \mathbf{e}_{i'} + \varepsilon \mathbf{e}_{i''} \succ \mathbf{x}, \quad 0 < \varepsilon < x_{i'} - x_{i''} \quad (9)$$

whenever $x_{i'} > x_{i''}$. The principle of transfers states that a transfer of any small amount from an outcome to any

other relatively worse–off outcome results in a more preferred outcome vector.

The rational preference relations satisfying additionally axioms (8) and (9) are called hereafter *fair (equitable) rational preference relations*. We say that outcome vector \mathbf{x}' *fairly dominates* \mathbf{x}'' ($\mathbf{x}' \succ_e \mathbf{x}''$), iff $\mathbf{x}' \succ \mathbf{x}''$ for all fair rational preference relations \succeq . In other words, \mathbf{x}' fairly dominates \mathbf{x}'' , if there exists a finite sequence of vectors \mathbf{x}^j ($j = 1, 2, \dots, s$) such that $\mathbf{x}^1 = \mathbf{x}''$, $\mathbf{x}^s = \mathbf{x}'$ and \mathbf{x}^j is constructed from \mathbf{x}^{j-1} by application of either permutation of coordinates, equitable transfer, or increase of a coordinate. Fig. 1 presents the structure of fair dominance for two-dimensional outcome vectors. For any outcome vector \mathbf{x} , the fair dominance relation distinguishes set $D(\mathbf{x})$ of dominated outcomes (obviously worse for all fair rational preferences) and set $S(\mathbf{x})$ of dominating outcomes (obviously better for all fair rational preferences). However, some outcome vectors are left (in white areas) and they can be differently classified by various specific fair rational preferences. The MMF assigns the entire interior of the inner white triangle to the set of preferred outcomes while classifying the interior of the external open triangles as worse outcomes. Isolines of various utility functions split the white areas in different ways. For instance, there is no fair dominance between vectors (0.01, 1) and (0.02, 0.02) and the MMF considers the latter as better while the so-called proportional fairness (PF) defined with logarithmic utility function [5] points out the former. On the other hand, vector (0.02, 0.99) fairly dominates (0.01, 1) and all fairness models (including MMF and PF) prefers the former.

An allocation pattern $\xi \in A$ is called *fairly (equitably) efficient* if $\mathbf{x} = \mathbf{f}(\xi)$ is fairly nondominated. Note that each fairly efficient solution is also Pareto-optimal, but not vice versa. The theory of majorization [10] includes the results which allow us to express the relation of fair (equitable) dominance

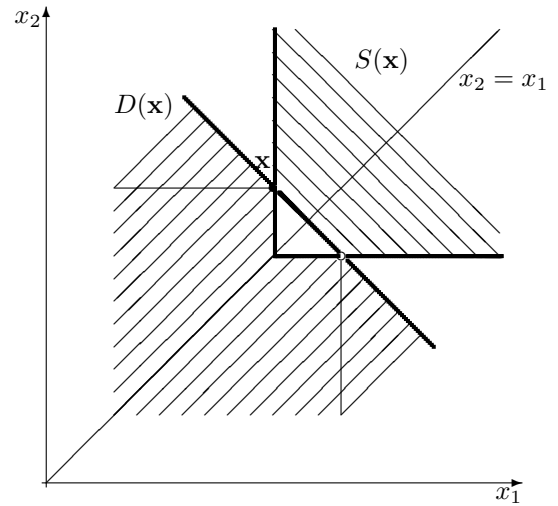


Fig. 1. Structure of the fair dominance: $D(\mathbf{x})$ – the set fairly dominated by \mathbf{x} . $S(\mathbf{x})$ – the set of outcomes fairly dominating \mathbf{x} .

as a vector inequality on the cumulative ordered outcomes [6]. This can be mathematically formalized as follows. First, introduce the ordering map $\Theta : R^m \rightarrow R^m$ such that $\Theta(\mathbf{x}) = (\theta_1(\mathbf{x}), \theta_2(\mathbf{x}), \dots, \theta_m(\mathbf{x}))$, where $\theta_1(\mathbf{x}) \leq \theta_2(\mathbf{x}) \leq \dots \leq \theta_m(\mathbf{x})$ and there exists a permutation τ of set I such that $\theta_i(\mathbf{x}) = x_{\tau(i)}$ for $i = 1, \dots, m$. Next, apply to ordered outcomes $\Theta(\mathbf{x})$, a linear cumulative map thus resulting in the *cumulative ordering map* $\bar{\Theta}(\mathbf{x}) = (\bar{\theta}_1(\mathbf{x}), \bar{\theta}_2(\mathbf{x}), \dots, \bar{\theta}_m(\mathbf{x}))$ defined as

$$\bar{\theta}_i(\mathbf{x}) = \sum_{j=1}^i \theta_j(\mathbf{x}), \quad i = 1, \dots, m \quad (10)$$

The coefficients of vector $\bar{\Theta}(\mathbf{x})$ express, respectively: the smallest outcome, the total of the two smallest outcomes, the total of the three smallest outcomes, etc. The theory of majorization allow us to derive the following theorem [6].

Theorem 1: Outcome vector \mathbf{x}' fairly dominates \mathbf{x}'' , if and only if $\bar{\theta}_i(\mathbf{x}') \geq \bar{\theta}_i(\mathbf{x}'')$ for all $i \in I$ where at least one strict inequality holds.

Theorem 1 permits one to express fair solutions of problem (7) as Pareto-optimal solutions to the multiple criteria problem

$$\max \{(\bar{\theta}_1(\mathbf{x}), \bar{\theta}_2(\mathbf{x}), \dots, \bar{\theta}_m(\mathbf{x})) : \mathbf{x} \in Q\}. \quad (11)$$

Indeed, the multiple criteria problem (11) may serve as a source of various fair allocation schemes [13], [15]. Although defined with simple linear constraints, the quantities $\bar{\theta}_k(\mathbf{x})$, used as criteria in (11), introduce m^2 additional variables and inequalities while m corresponds to the number of ordered pairs of network nodes which is already on the order of the square of the number of nodes $|V|$.

The ordered achievement vectors describe a distribution of outcomes generated by a given decision \mathbf{x} . In the case when there exists a finite set of all possible outcomes of the individual objective functions, we can directly deal with the distribution of outcomes described by frequencies of several outcomes. Let $V = \{v_1, v_2, \dots, v_r\}$ (where $v_1 < v_2 < \dots < v_r$) denote the set of all attainable outcomes (all possible values of the individual flows $x_i = f_i(\xi)$ for $\xi \in A$). We introduce integer functions $h_k(\mathbf{x})$ ($k = 1, \dots, r$) expressing the number of values v_k taken in the outcome vector \mathbf{x} . Having defined the functions h_k we can introduce cumulative distribution functions:

$$\bar{h}_k(\mathbf{x}) = \sum_{l=1}^k h_l(\mathbf{x}), \quad k = 1, \dots, r. \quad (12)$$

The function \bar{h}_k expresses the number of outcomes smaller or equal to v_k . Since we want to maximize all the outcomes, we are interested in the minimization of all the functions \bar{h}_k . The following assertion is valid [12]. For outcome vectors $\mathbf{x}', \mathbf{x}'' \in V^m$,

$$\Theta(\mathbf{x}') \geq \Theta(\mathbf{x}'') \Leftrightarrow \bar{\mathbf{h}}(\mathbf{x}') \leq \bar{\mathbf{h}}(\mathbf{x}''). \quad (13)$$

Note that $\bar{h}_r(\mathbf{x}) = m$ for any \mathbf{x} which means that the r -th quantity is always constant and therefore redundant.

In order to take into account the principle of transfers we need to distinguish values of outcomes smaller or equal to v_k . For this purpose we weight vector $\bar{\mathbf{h}}(\mathbf{x})$ to get:

$$\hat{h}_k(\mathbf{x}) = \sum_{l=1}^{k-1} (v_{l+1} - v_l) \bar{h}_l(\mathbf{x}) = \sum_{l=1}^{k-1} (v_k - v_l) h_l(\mathbf{x}) \quad (14)$$

for $k = 2, \dots, r$ and $\hat{h}_1(\mathbf{x}) = 0$. In other words, $\hat{h}_k(\mathbf{x})$ expresses the total of differences between v_k and all the outcomes x_i smaller than v_k . Since $(v_k - v_l) > 0$ for $1 \leq l < k$, it follows from (14) that vector function $\hat{\mathbf{h}}(\mathbf{x})$ provides a unique description of the distribution of coefficients of vector \mathbf{x} , i.e., for any $\mathbf{x}', \mathbf{x}'' \in V^m$ one gets:

$$\hat{\mathbf{h}}(\mathbf{x}') = \hat{\mathbf{h}}(\mathbf{x}'') \Leftrightarrow \mathbf{h}(\mathbf{x}') = \mathbf{h}(\mathbf{x}'') \Leftrightarrow \Theta(\mathbf{x}') = \Theta(\mathbf{x}'').$$

Moreover the following assertion is valid (Ogryczak, 1997). For achievement vectors $\mathbf{x}', \mathbf{x}'' \in V^m$,

$$\hat{\mathbf{h}}(\mathbf{x}') \leq \hat{\mathbf{h}}(\mathbf{x}'') \Leftrightarrow \bar{\Theta}(\mathbf{x}') \geq \bar{\Theta}(\mathbf{x}''). \quad (15)$$

Equivalence (15) permits one to express fair efficiency for problem (7) in terms of the standard efficiency for the multiple criteria problem with objectives $\hat{\mathbf{h}}(\mathbf{x})$:

$$\min \{(\hat{h}_1(\mathbf{x}), \hat{h}_2(\mathbf{x}), \dots, \hat{h}_r(\mathbf{x})) : \mathbf{x} \in Q\}. \quad (16)$$

Theorem 2: A feasible solution $\mathbf{x} \in Q$ is an fairly efficient solution of the multiple criteria problem (7), if and only if it is an efficient solution of the multiple criteria problem (16).

Formula (14) allows us to express $\hat{h}_k(\mathbf{x})$ as a piecewise linear function of \mathbf{x} :

$$\hat{h}_k(\mathbf{x}) = \sum_{i=1}^m \max\{v_k - x_i, 0\}, \quad k = 1, \dots, r. \quad (17)$$

Note that $\hat{h}_1(\mathbf{x}) = 0$ for any \mathbf{x} which means that the first criterion is constant and redundant in problem (16). Moreover, $mv_r - \hat{h}_r(\mathbf{x}) = \sum_{i=1}^m x_i$, thus representing the total throughput. Similarly, one may define for all k the complementary quantities $\eta_k(\mathbf{x}) = mv_k - \hat{h}_k(\mathbf{x}) = \sum_{i=1}^m \min\{x_i, v_k\}$ expressing the corresponding partial throughputs generated by flows ranged to v_k . Therefore, the entire multiple criteria model (16) can be reformulated as follows:

$$\begin{aligned} & \max [\eta_2, \eta_3, \dots, \eta_r] \\ & \text{s.t.} \\ & \eta_k = \sum_{i=1}^m t_{ki}, \quad k = 2, \dots, r, \\ & t_{ki} \leq x_i, \quad i = 1, \dots, m; \quad k = 2, \dots, r, \\ & t_{ki} \leq v_k, \quad i = 1, \dots, m; \quad k = 2, \dots, r, \\ & \mathbf{x} \in Q \end{aligned} \quad (18)$$

Note that the above formulation adds only linear constraints to the original attainable set Q . Hence, for the basic network dimensioning problems with the set Q defined by constraints (1)–(6), the resulting formulation (18) remains in the class of (multiple criteria) MILP.

IV. MULTIPLE TARGET ANALYSIS

Although defined with simple linear constraints, the expanded model (18) introduces $r \times m$ additional variables and inequalities. This may cause a serious computational burden for real-life network dimensioning problems. Note that the number of services (traffic demands) m corresponds to the number of ordered pairs of network nodes which is already square of the number of nodes $|V|$. On the other hand, quantity r represents the number of various possible outcomes (flow sizes). In order to reduce the problem size one may attempt to restrict the number of distinguished outcome values (criteria in the problem (18)).

Let us consider a sequence of indices $K = \{k_1, k_2, \dots, k_q\}$, where $v_{k_1} < v_{k_2} < \dots < v_{k_{q-1}} < v_{k_q}$, and the corresponding restricted form of the multiple criteria model (16):

$$\max \{(\eta_{k_1}, \dots, \eta_{k_q}) : \mathbf{x} \in Q\} \quad (19)$$

with only $q < r$ criteria. Following Theorem 2, multiple criteria model (16) allows us to generate any fairly efficient solution of problem (7). Reducing the number of criteria we restrict these opportunities. Nevertheless, one may still generate reasonable compromise solutions. First of all the following assertion is valid.

Theorem 3: If \mathbf{x}^o is an efficient solution of the restricted problem (19), then it is an efficient (Pareto-optimal) solution of the multiple criteria problem (7) and it can be fairly dominated only by another efficient solution \mathbf{x}' of (19) with exactly the same values of criteria: $\hat{h}_k(\mathbf{x}') = \hat{h}_k(\mathbf{x}^o)$ for all $k \in K$.

Proof: Suppose, there exists $\mathbf{x}' \in Q$ which dominates \mathbf{x}^o . This means, $x'_i \geq x_i^o$ for all $i \in I$ with at least one inequality strict. Hence, $\hat{h}_k(\mathbf{x}') \leq \hat{h}_k(\mathbf{x}^o)$ for all $k = 1, \dots, r$ and $\hat{h}_{k_q}(\mathbf{x}') < \hat{h}_{k_q}(\mathbf{x}^o)$ which contradicts efficiency of \mathbf{x}^o within the restricted problem (19).

Suppose now that $\mathbf{x}' \in Q$ fairly dominates \mathbf{x}^o . Due to Theorem 2, this means that $\hat{h}_k(\mathbf{x}') \leq \hat{h}_k(\mathbf{x}^o)$ for all $k = 1, \dots, r$ with at least one inequality strict. Hence, $\hat{h}_k(\mathbf{x}') \leq \hat{h}_k(\mathbf{x}^o)$ for all $k \in K$ and any strict inequality would contradict efficiency of \mathbf{x}^o within the restricted problem (19). Thus, $\hat{h}_k(\mathbf{x}') = \hat{h}_k(\mathbf{x}^o)$ for all $k \in K$. ■

It follows from Theorem 3 that while restricting the number of criteria in the multiple criteria model (16) we can essentially still expect reasonably fair efficient solution and only *unfairness* may be related to the distribution of flows within classes of skipped criteria. In other words, we have guaranteed some rough fairness while it can be possibly improved by redistribution of flows within the intervals $(v_{k_j}, v_{k_{j+1}}]$ for $j = 1, 2, \dots, q-1$. Since the fairness preferences assume increase of smaller flows against larger ones, they aware of the use of very small flows. One may introduce a grid of critical values $v_{k_1} < v_{k_2} < \dots < v_{k_{q-1}} < v_{k_q}$ which is dense for smaller indices (smaller flow values) while sparser for larger indices thus expecting some reasonable compromise between fairness and throughput maximization. In our computational analysis on the network with 132 elastic flows and the total throughput requirements ranging between 500 and 1100 (Section V) we have preselected 11 values v_k as 1, 2, ..., 10, and 20.

Finally, we may generate various fairly efficient bandwidth allocation patterns as efficient solutions of the multiple criteria problem:

$$\begin{aligned} \max \quad & (\eta_k)_{k \in K} \\ \text{s.t.} \quad & \mathbf{x} \in Q \\ & \eta_k = \sum_{i \in I} t_{ki}, \quad k \in K \\ & t_{ki} \leq x_i, \quad t_{ki} \leq v_k, \quad i \in I, k \in K \end{aligned} \quad (20)$$

where the attainable set Q is defined by constraints (1)–(6). Exactly, in the case of the complete multiple criteria model ($K = \{1, \dots, r\}$), according to Theorem 2, all fairly efficient allocations can be found as efficient solutions to (20) while in the case of restricted set of criteria K some minor unfairness related to the distribution of flows within classes of skipped criteria may occur (Theorem 3).

The simplest way to model a large gamut of fairly efficient allocations may depend on the use some combinations of criteria $(\eta_k)_{k \in K}$. Better controllability and the complete parameterization of nondominated solutions for discrete problems can be achieved with the direct use of the reference point methodology introduced by Wierzbicki [19] and later extended leading to efficient implementations of the so-called aspiration/reservation based decision support (ARBDS) approach [8]. The ARBDS approach allows the decision maker (DM) to specify the requirements in terms of aspiration and reservation levels, i.e., by introducing acceptable and required values for several criteria. Depending on the specified aspiration and reservation levels, a special scalarizing achievement function is built and maximized. Maximization of the scalarizing achievement function generates an efficient solution to the multiple criteria problem. The solution is accepted by the DM or some modifications of the aspiration and reservation levels are introduced to continue the search for a better solution. When applying the ARBDS methodology to the multiple target model (20), one may generate various fairly efficient solutions of the original problem (7).

While building the scalarizing achievement function the following properties of the preference model are assumed. First of all, the function must be strictly increasing with respect to each outcome to guarantee that more is preferred to less (maximization) for any individual outcome η_k . Second, a solution with all individual outcomes η_k satisfying the corresponding reservation levels is preferred to any solution with at least one individual outcome worse (smaller) than its reservation level. Next, provided that all the reservation levels are satisfied, a solution with all individual outcomes η_k equal to the corresponding aspiration levels is preferred to any solution with at least one individual outcome worse (smaller) than its aspiration level. That means, the scalarizing achievement function maximization must enforce reaching the reservation levels prior to further improving of criteria. In other words, the reservation levels represent some soft lower bounds on the maximized criteria. When all these lower bounds are satisfied, then the optimization process attempts to reach the aspiration levels.

The generic scalarizing achievement function takes the

following form [19]:

$$\sigma(\eta) = \min_{k \in K} \{\sigma_k(\eta_k)\} + \varepsilon \sum_{k \in K} \sigma_k(\eta_k) \quad (21)$$

where ε is an arbitrary small positive number and σ_k , for $k \in K$, are the partial achievement functions measuring actual achievement of the individual outcome η_k with respect to the corresponding aspiration and reservation levels (η_k^a and η_k^r , respectively). Thus the scalarizing achievement function is, essentially, defined by the worst partial (individual) achievement but additionally regularized with the sum of all partial achievements. The regularization term is introduced only to guarantee the solution efficiency in the case when the maximization of the main term (the worst partial achievement) results in a non-unique optimal solution.

The partial achievement function σ_k can be interpreted as a measure of the DM's satisfaction with the current value (outcome) of the k -th criterion. It is a strictly increasing function of outcome η_k with value $\sigma_k = 1$ if $\eta_k = \eta_k^a$, and $\sigma_k = 0$ for $\eta_k = \eta_k^r$. Thus the partial achievement functions map the outcomes values onto a normalized scale of the DM's satisfaction. We use the piecewise linear partial achievement function introduced in [12]:

$$\sigma_k(\eta_k) = \begin{cases} \gamma(\eta_k - \eta_k^r)/(\eta_k^a - \eta_k^r), & \eta_k \leq \eta_k^r \\ (\eta_k - \eta_k^r)/(\eta_k^a - \eta_k^r), & \eta_k^r < \eta_k < \eta_k^a \\ \beta(\eta_k - \eta_k^a)/(\eta_k^a - \eta_k^r) + 1, & \eta_k \geq \eta_k^a \end{cases}$$

where β and γ are arbitrarily defined parameters satisfying $0 < \beta < 1 < \gamma$. In our implementation the values $\beta = 0.01$ and $\gamma = 100$ have been used. This partial achievement function is strictly increasing and concave which guarantees its LP computability with respect to outcomes η_k .

Recall that in our model outcomes η_k represent partial throughputs for ranged flows x_i , i.e. $\eta_k = \sum_{i=1}^m \min\{x_i, v_k\}$. Hence, the reference vectors (aspiration and reservation) represent, in fact, some reference distributions of outcomes (flows). Moreover, due to the cumulation of outcomes, while considering equal flows ϕ as the reference (aspiration or reservation) distribution, one needs to set the corresponding levels as $\eta_k = m\phi$ for $\phi \leq v_k$ and $\eta_k = mv_k$ otherwise. Certainly, one may specify any desired reference distribution in terms of ranged throughputs. Although, special meaning of the last (throughput) criterion should be rather operated independently from the others. Such an approach to control the search for a compromise fair and efficient bandwidth allocation has been confirmed by the computational experiments as described in the following section.

V. COMPUTATIONAL EXAMPLES

The reference distribution approach described in preceding sections has been tested on a sample network dimensioning problem with elastic traffic. The network topology of the presented problem (Fig. 2) is patterned after the backbone network of a Polish ISP [13]. The network consists of 12 nodes and 18 links. Flows between any pair of different nodes have been considered (i.e., $144 - 12 = 132$ flows). For each flow,

two alternative paths have been specified that could be used for transport. All information of a flow had to travel along one of the paths. All links have unit costs equal to one, and the budget for link bandwidth is $B = 1000$. Since all links have equal costs of one, path cost are equal to the path length (1, 2, 3 or 4 for the shortest paths in the example topology). For each flow, two alternative paths (the shortest and the second shortest) have been specified that could be used for transport. The entire flow had to travel along one of the paths with no splitting allowed (nonbifurcation formulation (1)–(2)).

We have analyzed the network dimensioning problem defined by constraints (1)–(6). Thus the model under consideration allows flows to choose one of two paths for transport (1)–(2) and limits the capacity of certain links from above while providing also some free link capacity for certain links (3). The intention behind the model has been to depict a situation when the network operator wishes to extend the capacity of an existing network. In this network, certain links cannot be upgraded beyond a certain values to the use of legacy technologies, due to prohibitive costs or administrative reasons (for instance, it may be cheap to use already installed fiber that has not been in use before, but it may be prohibitively expensive to install additional fiber). Actually, free link capacity was set to 10, and the upper limit on the expansion capacity was set to 30. Due to the presence of free link capacity and upper limits on link capacity, the MILP solver found solutions where certain flows had to use alternative paths rather than the shortest paths. These flows were more expensive than other flows that were allowed to use their shortest paths. Recall that we have used a single-path formulation, meaning that the entire flow had to be switched to the alternative path. Flows could not be split, which is consistent with several traffic engineering technologies used today.

For all model versions, the final input to the model consisted of the reservation and aspiration levels for the total throughput within ranges of the specified target flow values. We have preselected 11 target flow values: ten values $v_k = k$ ($k = 1, 2, \dots, 10$), and $v_{11} = 20$. For simplicity of the analysis, all aspiration levels were set larger than the maximum possible value $mv_k + 1$, and only reservation levels were used to control the outcome flows. One of the most significant parameters was the reservation level for the largest target value representing

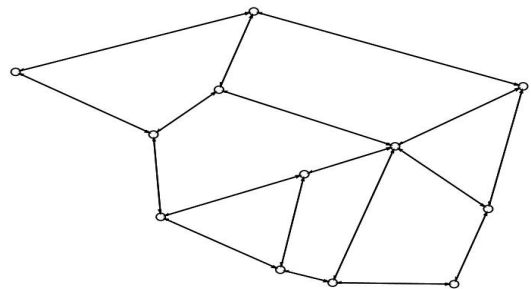


Fig. 2. Sample network topology patterned after the backbone network of a Polish ISP.

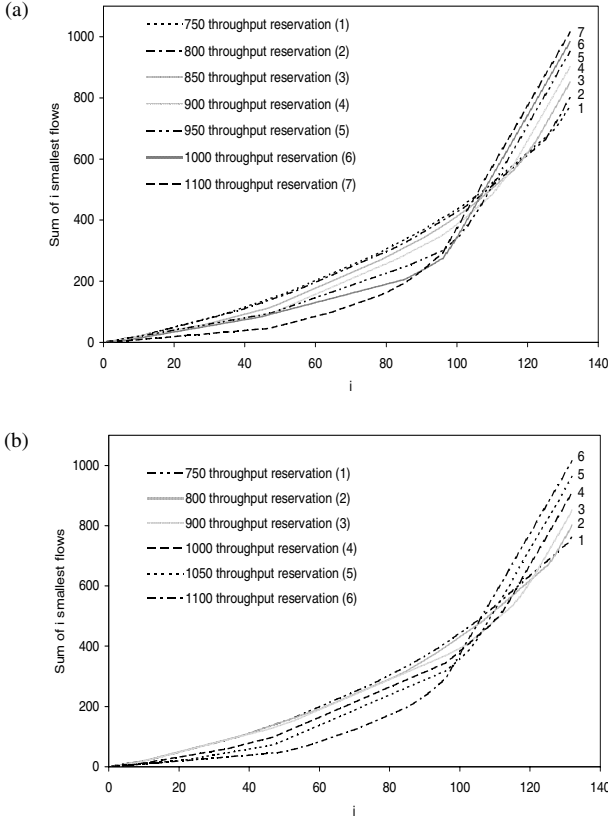


Fig. 3. Flow distribution for varying throughput reservation with $\phi = 2$ and $s = 0.01$ (a), $s = 0.04$ (b).

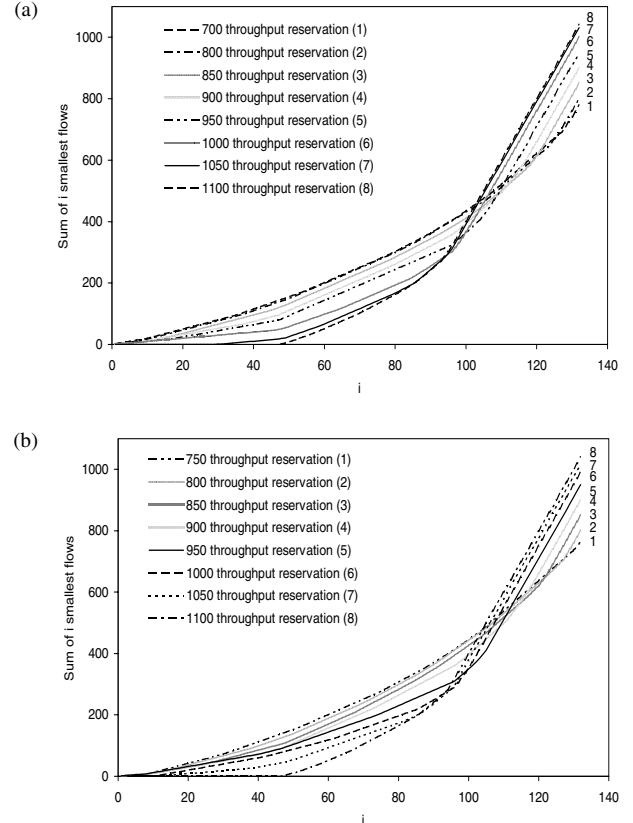


Fig. 4. Flow distribution for varying throughput reservation with $\phi = 1$ and $s = 0.01$ (a), $s = 0.04$ (b).

actually the required network throughput. This value denoted by η_r^r was selected separately from the other reservation levels. All the other reservation levels were formed as a linearly increasing sequence of the ordered values with slope (step) s . Exactly, a value ϕ is selected as required minimal flow and further the (ordered) required flows are defined as $\phi + (i-1)s$. Hence, one gets the reservation levels $\eta_k^r = mv_k$ for $v_k \leq \phi$ and $mv_k - k(k-1)s/2$ where $\bar{k} = (v_k - \phi)/s$. For the sake of simplicity, we select value ϕ as one of the lower target values v_k . Thus we have 3 control parameters: reservation level for total throughput, minimum required flow, and the slope of ordered required flows.

Fig. 4 and 3 present plots of cumulated ordered flows $\bar{\theta}_i(\mathbf{x})$ versus number i (rank of a flow in ordering according to flow throughput). The total network throughput is represented in the figures by the altitude of the right end of the curve ($\bar{\theta}_{132}(\mathbf{x})$). A perfectly equal distributions of flows would be graphically represented by an ascending line of constant slope. Solutions resulting in similar flow distribution have been skipped.

In the first experiment, we set the minimal flow value $\phi = 1$ while the throughput reservation level η_m^r and the slope s have been used to search for compromise solutions that traded off fairness against efficiency. The throughput reservation level has been varied from 500 to 1100 for two values of slope $s = 0.01$ and 0.04 . For low throughput reservations (up to 750) no significant conflict with fairness has occurred and the

total throughput about 750 has been provided with up to 10 flows on level of 2 and other larger than 2 (for $s = 0.01$) and up to 8 flows on the level of 1 and other larger than 1 (for $s = 0.04$). As η_m^r increases, the cheaper flows receive more throughput at the expense of more expensive (longer) flows. For $s = 0.01$ and values of η_m^r above 1000, some flows are starved; almost 50 for $\eta_m^r = 1100$. Actually, no solution has managed to reach the total throughput significantly larger than 1050. For $s = 0.04$, already values of η_m^r above 950 cause that a few flows are starved. When repeating the experiment with the required minimal flow value $\phi = 2$, for $s = 0.01$ and values of η_m^r up to 950 all flows reach the level of at least 2. For larger values of η_m^r up to 1100 still all flows remain positive but a few of them are reduced to the level of 1. Finally, for throughput reservation 1100 about 1% of flows are starved. For $s = 0.04$ and values of η_m^r up to 800 all flows reach the level of at least 2. For larger values, all flows remain positive.

Overall, the experiments on the sample network topology demonstrated the versatility of the described methodology for fair network dimensioning. The use of reservation levels, controlled by a small number of simple parameters, allowed us to search for solutions best fitted to various possible preferences of a network designer. Moreover, the computation time to generate a single solution of our sample dimensioning problem does not exceed 200 seconds when using CPLEX 9.1 on an 1.4GHz PC, thus enabling an interactive search of

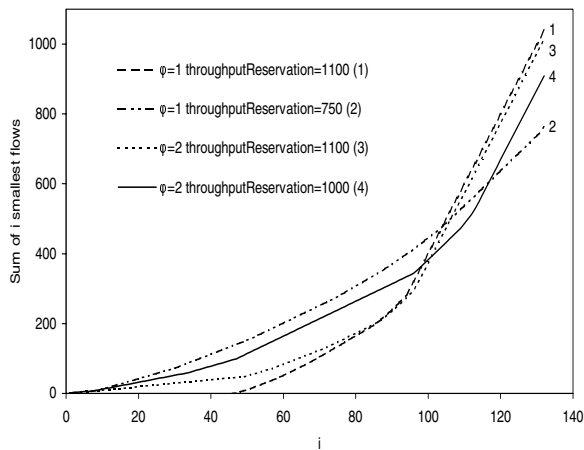


Fig. 5. Flows distribution for varying parameters in the interactive analysis.

a satisfactory fair and efficient allocation. Consider a network designer who wishes to extend link capacities of a network shown in Fig. 2. Note that the network designer need not have a set of directly expressed preferences, but rather is an expert that works using tacit knowledge. First, she may search for a solution that has a high overall throughput by setting 1100 as the corresponding reservation, and some reasonable slope value $s = 0.04$ and the required minimal flow value $\phi = 1$. Solution 1, shown in Fig. 5, indeed provides a high throughput (about 1050) but is unfair with more than 40 flows completely starved. When the network designer relaxes the throughput reservation to 750, she gets Solution 2 that is quite close to perfectly fair, but has a low throughput (750). To identify a relatively fair solution with a larger throughput, the network designer returns to $\eta_m^r = 1100$, but increases the required smallest flow to $\phi = 2$. This leads to Solution 3 with all positive flows, (but about 40 of them below 2) and a total throughput of about 1000. The solution plot is very similar to that of Solution 1, thus depicting strong unfairness. To reach a more fair solution, the network designer accepts a decrease of the total throughput. When setting $\eta_m^r = 1000$, she gets Solution 4 with a total throughput about 900, but similarly fair as Solution 2. In Solution 4, about 80% of (smallest) flows remain on the level close to those of Solution 2, while about 20% receive much larger values, thus increasing the total throughput. The network designer finds Solution 4 an acceptable compromise. The presented analysis is simplified, but it demonstrates that it is possible to easily find a satisfactory fair and efficient allocation pattern in a few interactive steps. Moreover, the plots of cumulated ordered flows turn out to be a convenient graphical interface to support the search process. The selected solution contains values for link capacities and allows the network designer to extend the network in Fig. 2.

VI. CONCLUSION

A central issue in networking is how to allocate bandwidth to flows efficiently and fairly. The Max-Min Fairness is widely used to meet these goals. Allocating the resources to optimize

the worst performances may cause a large worsening of the overall (mean) performances. Therefore, several other fair allocation schemes have been researched and analyzed. We have shown that there exists a multiple criteria model that allows to represent consistently the overall efficiency and fairness goals. Moreover, the criteria measure actual network throughputs for various levels (targets) of flows. Thereby, the criteria can easily be introduced into the model. While looking for fairly efficient bandwidth allocation the reference point methodology can be applied to the multiple target partial throughputs. Our initial experiments with such an approach to the problem of network dimensioning with elastic traffic have confirmed the theoretical advantages of the method. We were easily able to generate various (compromise) fair solutions despite the fact that the search for fairly efficient compromise solutions was controlled by only three parameters. One of these parameters was a reservation level for the network throughput.

ACKNOWLEDGMENT

The research was supported by the Ministry of Science and Information Society Technologies under grant 3T11C 005 27 (Włodzimierz Ogryczak and Adam Wierzbicki) and under grant 3T11D 001 27 (Marcin Milewski).

REFERENCES

- [1] Bertsekas D, Gallager R. Data Networks. Englewood Cliffs: Prentice-Hall, 1987.
- [2] Bonald T, Massoulié L. Impact of fairness on Internet performance. Proceedings of ACM Sigmetrics, June 2001, 82–91.
- [3] Denda R, Banachs A, Effelsberg W. The fairness challenge in computer networks. Lect Notes in Comp Sci 2000; 1922: 208–220.
- [4] Jaffe J. Bottleneck flow control. IEEE Trans. on Communications 1980; 7:207–237.
- [5] Kelly F, Mauloo A, Tan D. Rate control for communication networks: shadow prices, proportional fairness and stability. J Oper Res Soc 1997; 49:206–217.
- [6] Kostreva MM, Ogryczak W. Linear optimization with multiple equitable criteria. RAIRO Oper Res 1999; 33:275–297.
- [7] Kostreva MM, Ogryczak W, Wierzbicki A. Equitable aggregations and multiple criteria analysis, European J Opnl Res 2004; 158:362–367.
- [8] Lewandowski A, Wierzbicki AP. Aspiration Based Decision Support Systems — Theory, Software and Applications. Berlin: Springer, 1989.
- [9] Luss H. On equitable resource allocation problems: a lexicographic minimax approach. Oper Res 1999; 47:361–378.
- [10] Marshall AW, Olkin I. Inequalities: Theory of Majorization and Its Applications. New York: Academic Press, 1979.
- [11] Mo J, Walrand J. Fair end-to-end window-based congestion control. IEEE/ACM Trans on Networking 2000; 8:556–567.
- [12] Ogryczak W. Linear and Discrete Optimization with Multiple Criteria: Preference Models and Applications to Decision Support (in Polish). Warsaw: Warsaw Univ Press, 1997.
- [13] Ogryczak W, Śliwiński T, Wierzbicki A. Fair resource allocation schemes and network dimensioning problems. J Telecomm and Info Tech 2003; 3:34–42.
- [14] Ogryczak W, Tamir A. Minimizing the sum of the k largest functions in linear time. Information Proc Letters 2003; 85:117–122.
- [15] Ogryczak W, Wierzbicki A. On multi-criteria approaches to bandwidth allocation. Control and Cybernetics 2004; 33:427–448.
- [16] Pióro M, Medhi D. Routing, Flow and Capacity Design in Communication and Computer Networks. San Francisco: Morgan-Kaufmann, 2004.
- [17] Rawls J. The Theory of Justice. Cambridge: Harvard Univ Press, 1971.
- [18] Tang A, Wang J, Low SH. Is fair allocation always inefficient. IEEE INFOCOM 2004.
- [19] Wierzbicki AP. A mathematical basis for satisficing decision making. Math Modelling 1982, 3:391–405.