

Distributed Fair Resource Allocation in Cellular Networks in the Presence of Heterogeneous Delays

Lei Ying, R. Srikant, Atilla Eryilmaz, and Geir E. Dullerud

Abstract—We consider the problem of allocating resources at a base station to many competing flows, where each flow is intended for a different receiver. The channel conditions may be time-varying and different for different receivers. It has been shown in a previous paper that in a delay-free network, a combination of queue-length-based scheduling at the base station and congestion control at the end users can guarantee queue-length stability and fair resource allocation. In this note, we extend this result to wireless networks where the congestion information from the base station is received with a feedback delay at the transmitters.

Index Terms—Congestion control, resource allocation, queue-length-based policy, cellular networks, heterogeneous delays.

I. INTRODUCTION

We study the problem of fair allocation of resources in the downlink of a cellular wireless network consisting of a single base station and many receivers. The data destined for each receiver is maintained in a separate buffer. The arrivals to the buffers are determined via a congestion control mechanism, which will be described in detail later. We assume that the time is slotted. The channels between the base station and the receivers are assumed to have random time-varying gains which are independent from one time-slot to the next. The independence assumption can be relaxed easily, but we use it here for ease of exposition. The goal is to allocate the network capacity fairly among the users, in accordance with the needs of the users, while exploiting the time-variations in the channel conditions. We associate a utility function with each user that is a concave, increasing function of the mean service that it receives from the network. In [6], it was shown that a combination of Internet style congestion control at the end-users and queue-length based scheduling at the base station achieves the goal of fair and stabilizing resource allocation. This result is somewhat surprising since the resource constraints in the case of a wireless network are very different from the linear constraints in the case of the Internet [16]. The relative merits of congestion control-based resource allocation scheme as compared to other resource allocation schemes for cellular networks are discussed in [6]. Several other works in the same context are [17], [10], and [13].

In [6], it is assumed that there are no delays in the transmission of packets from an end-user (transmitter) to the base station and in the transmission of congestion information from the base station back to the end users. However, if we consider the case where the end users are connected to the base station through the Internet, then delays exist in both directions: There is a propagation delay τ_i^f from the end user

i to the base station—we call it the forward delay of the end user i , and a propagation delay τ_i^b from the base station to the end user i —we call it the backward delay. It is well-known that the presence of delays may affect the performance of the network. For example, Internet congestion controllers which are globally stable for the delay-free network may become unstable if the feedback delays are large [16]. In our problem, when delays exist, the information the end users obtain will be “outdated” information. So the congestion information the users obtain at time t does not reflect the queue status at the base station at time t . So it is interesting to study a wireless network with delays and ask whether the conclusions of [6] still hold for wireless networks with heterogeneous delays. We answer this question by showing that for a network with uniformly-bounded delays, which are potentially heterogeneous and time-varying, the algorithm of [6] is stable and can be used to approximate weighted- m fair allocation arbitrarily closely. We emphasize that the results hold for networks with arbitrarily large, but bounded time-varying delays. So even if the end users can only get very old feedback information from the base station, the network is still stable and can approach the fair resource allocation. On the other hand, from the proof, we can also see that when the delays are large, it may take more time for the network to achieve the fair resource allocation. This observation is also supported by simulations, not shown here due to page limitations, which are presented in [20].

II. SYSTEM MODEL

We consider a cellular network shared by n flows in the downlink and assume that the base station maintains n separate queues, one corresponding to each flow. We study the fair resource allocation problem in this note. Specifically we consider weighted- m fairness. It means that each source i has a utility function given by $U_i(\bar{z}_i) = \alpha_i \frac{\bar{z}_i^{1-m}}{1-m}$, where \bar{z}_i is the average rate at which user i transmits and α_i is a positive weighting factor [12]. Here, $m = 1$, $m = 2$ and $m \rightarrow \infty$ correspond to respectively proportional fair, minimum potential delay fair and max-min fair resource allocations. Note that in the case, $m = 1$, the aforementioned utility function is not well-defined; however, it can be shown that as $m \rightarrow 1$, the limiting resource allocation corresponds to the case where $U_i(\bar{z}_i) = \alpha_i \log \bar{z}_i$ [12]. We assume that the time is slotted and denote the length of the queue i at the beginning of the time slot t by $x_i[t]$, the number of arrivals to queue i in time slot t by $a_i[t]$, and the amount of service offered to queue i in slot t by $\mu_i[t]$. We assume that each of these parameters can only take non-negative integer values. The evolution of the size of the i th queue is given by

$$x_i[t + 1] = x_i[t] + a_i[t] - \mu_i[t] + u_i[t]$$

where $u_i[t]$ is a non-negative quantity which denotes the wasted service given to queue i at time slot t and it guarantees that $x_i[t] \geq 0$. We also assume that the channel between the base station and the receivers can be in one of J states in a given slot. We use $s[t]$ to denote the state in time slot t . The channel state is assumed to be fixed within a time slot, but may vary from one slot to another, thus capturing the time-varying characteristics of a fading environment. Corresponding to each channel state, say j , is an achievable rate region, C_j , that is defined to be convex hull of the feasible rate vectors, $\eta[t] := (\eta_1[t], \dots, \eta_n[t])$, that can be offered to the queues. We assume that each C_j is a bounded region and let $\hat{\eta} < \infty$ denote the upper bound on the achievable rates for all channel states. The channel state process is assumed to be independent and identically distributed in each time slot, but we do not require that

Manuscript received May 16, 2005; revised November 7, 2005 and August 13, 2006. Recommended by Associate Editor L. Xie. This work was supported by the Air Force Office of Scientific Research under Grant F49620-01-1-0365, and by the National Science Foundation under Grant ECS-0401125.

L. Ying, R. Srikant, and G. E. Dullerud are with Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA (e-mail: lying@uiuc.edu; rsrikant@uiuc.edu; dullerud@uiuc.edu).

A. Eryilmaz is with the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: eryilmaz@mit.edu)

Digital Object Identifier 10.1109/TAC.2006.886531

the statistics be known at the base station. Furthermore, we define the mean achievable rate region as

$$\bar{C} := \left\{ \eta : \eta = \sum_{j=1}^J \pi_j^{\text{ch}} \eta^{(j)}, \eta^{(j)} \in C_j \right\}$$

where π_j^{ch} stands for the stationary distribution of the channel state process being in state j . The scheduler will use following algorithm.

SCHEDULER: Given the current queue length $\mathbf{x}[t] := (x_1[t], \dots, x_n[t])$ and current channel state $s[t]$, the scheduler at the base station chooses a service rate vector $\mu[t] := (\mu_1[t], \dots, \mu_n[t]) \in C_{s[t]}$ that satisfies

$$\mu[t] \in \arg \max_{\eta \in C_{s[t]}} \sum_{i=1}^n x_i[t] \eta_i.$$

This scheduling rule was introduced in the context of fixed arrival rates (i.e., where the arrival rates are not adjusted by a congestion controller) in [18], where it was also shown that it is a stabilizing rule, i.e., the mean queue-lengths are upper-bounded. This result was extended in many different directions in [2], [15], [7], [3], [9], [5], and [14].

In our model, the packet arrival rate into the queue is assumed to be controlled according to the well-known dual controller that has been studied extensively in the context of Internet congestion control [8], [11], [19], [16]. In the context of Internet congestion control, a dual controller chooses the transmission rate z_i such that

$$\frac{\alpha_i K}{x_i} = z_i$$

for any some $K > 0$. Next, we describe the operation of our congestion controller followed by some assumptions.

CONGESTION CONTROLLER: Recall that, for user i , the forward delay is τ_i^f and the backward delay is τ_i^b . In our model, the downlink is the only bottleneck of the system, so τ_i^f , which is the propagation delay from user i to the base station, is a constant. On the other hand, packets experience queuing delay at the base station, and the transmissions between the base station and the receivers are over wireless links. Thus, $\tau_i^b[t]$ is time varying. Since users will always use the latest feedback information, we define $\tau_i^b[t]$ such that

$$\tau_i^b[t] = \min\{\tilde{\tau}_i^b[t], \tau_i^b[t-1] + 1\}$$

where $\{\tilde{\tau}_i^b[t]\}_t$ are i.i.d. random variables, and $T_{\max} - \tau_i^f \geq \tilde{\tau}_i^b[t] \geq \tau_i^p$. Note that T_{\max} is the upper bound on the round trip delays, and τ_i^p is the propagation delay from the base station to user i via receiver i . Now, denote the amount of data sent out by user i in slot t by $\lambda_i[t]$. The congestion controller at user i regulates the mean of $\lambda_i[t]$ such that

$$E \left[\lambda_i[t] \mid x_i \left[t - \tau_i^b[t] \right] \right] = \min \left\{ \frac{\alpha_i K}{(x_i[t - \tau_i^b[t]])^m}, M \right\} \quad (1)$$

where $m > 0, M > 2\hat{\eta}$ is a positive constant which ensures that the arrival rate into the queue is upper bounded when the queue length is small, and $x_i[t - \tau_i^b[t]]$ is the congestion information measured by the based station and feedback to user i via receiver i . We will later

show that K has to be large to approximate weighted- m fair resource allocation.

Since $a_i[t] = \lambda_i[t - \tau_i^f]$, the mean of the number of arrivals into queue i at time t given by

$$E[a_i[t] | x_i[t - T_i[t]]] = E \left[\lambda_i \left[t - \tau_i^f \right] \mid x_i \left[t - \tau_i^f - \tau_i^b \left[t - \tau_i^f \right] \right] \right]$$

thus we have

$$E[a_i[t] | x_i[t - T_i[t]]] = \min \left\{ \frac{\alpha_i K}{(x_i[t - T_i[t]])^m}, M \right\} \quad (2)$$

where $T_i[t] = \tau_i^f + \tau_i^b[t - \tau_i^f]$. Define $\tilde{T}_i[t] = \tau_i^f + \tilde{\tau}_i^b[t - \tau_i^f]$, we have $T_i[t] = \min\{\tilde{T}_i[t], T_i[t-1] + 1\}$, and $\{\tilde{T}_i[t]\}_t$ are i.i.d. random variables such that $T_{\max} \geq \tilde{T}_i[t] \geq \tau_i^f + \tau_i^p$. We assume $a_i[t]$ is independent across time slots and

$$E[a_i^2[t] | x_i[t - T_i(t)]] \leq V < \infty \text{ for all } x_i[t - T_i(t)]. \quad (3)$$

Furthermore, we assume there exist positive numbers $\theta, A > T\hat{\eta}$, and $h > 2$ such that for any $N > A$

$$P \left(\sum_{j=1}^{T_{\max}} a_i[t-j] = N \right) < \frac{\theta}{N^h}, \text{ for all } i \text{ and all } t. \quad (4)$$

In summary, the combined scheduler-congestion controller algorithm can be defined as follows:

$$x_i[t+1] = x_i[t] + a_i[t] - \mu_i[t] + u_i[t] \quad (5)$$

$$\mu[t] \in \arg \max_{\eta \in C_{s[t]}} \sum_{i=1}^n x_i[t] \eta_i \quad (6)$$

where $a_i[t]$ is a random variable satisfying the conditions in (2)–(4). Note that the congestion control part of this algorithm is slightly different from the algorithm in [6]. We impose an upper-bound on the source rates in a more natural manner than in [6]. Our results continue to hold for the algorithm in [6] too.

We now present the following theorem, which will be useful later. This theorem is similar to [6, Prop. 1].

Theorem 1: There exists a unique pair of vectors (\mathbf{x}^*, μ^*) which satisfy the following conditions:

- $\mu^* \in \arg \max_{\eta \in \bar{C}} \sum_{i=1}^n x_i^* \eta_i$;
- $x_i^* = \left(\frac{\alpha_i K}{\mu_i^*} \right)^{(1/m)}$ for all i ;
- μ^* is the optimal solution to $\max_{\mu \in \bar{C}} \sum_{i=1}^n K \alpha_i \frac{\mu_i^{1-m}}{1-m}$.

From the previous theorem, we can see that μ^* is weighted- m fair. For the stochastic model, we will show that $\mu[t]$ converges to μ^* , defined in Theorem 1, in a probabilistic sense. This then implies that the network reaches a fair operating point.

III. WEIGHTED- m FAIRNESS AND STABILITY

Use $\mu[\infty]$ to denote the steady state of μ , we will show in our main theorem—Theorem 3, that for any $\epsilon < 0$

$$\lim_{K \rightarrow \infty} P(|\mu[\infty] - \mu^*| \geq \epsilon) = 0$$

which implies that the network can approximate the weighted- m fairness when K is chosen to be large. To prove this, we first need following lemma which characterizes the mean distance between \mathbf{x}^* and the steady state of $\mathbf{x}[t]$.

Lemma 2: There exists a positive constant $\sigma < 1/m$, and a positive constant \bar{c} that depends on the mean achievable rate region, the algorithm parameters $\{\alpha_i\}$, and the moments of the channel and arrival process, such that

$$E[\|\mathbf{x}[\infty] - \mathbf{x}^*\|] \leq \bar{c}K^{\frac{1}{m}-\sigma} \text{ for large } K$$

where $\mathbf{x}[\infty]$ is an informal notation for the steady state of \mathbf{x} and $\|\cdot\|$ denotes the Euclidean distance in the \mathfrak{R}^n .

Proof: Define $\mathbf{y}[t] = (\mathbf{x}[t], \dots, \mathbf{x}[t - T_{\max}], \mathbf{T}[t])$, where $T_{\max} \geq T_i[t]$. It is easy to see that the process $\{\mathbf{y}[t]\}_{t \geq 0}$ is a Markov chain because $a_i[t]$ depends only on $x_i[t - T_i[t]]$ and $T_i[t]$ depends only on $T_i[t - 1]$, so $x_i[t + 1]$, $T_i[t + 1]$ and $\mathbf{y}[t + 1]$ are determined by $\mathbf{y}[t]$. Further, define the Foster–Lyapunov function

$$W(\mathbf{y}[t]) = \frac{1}{2} \sum_i^n (x_i[t] - x_i^*)^2$$

and

$$E[\Delta W_t(\mathbf{y})] := E[W(\mathbf{y}[t + 1]) - W(\mathbf{y}[t]) | \mathbf{y}[t]].$$

Then following an argument similar to [6, Th. 2], the lemma will hold if there exist a finite set S_σ , positive numbers $\sigma < 1/m$, δ^* , and ζ such that for large K

$$E[\Delta W_t(\mathbf{y})] \leq -\frac{\delta^*}{K^{\frac{1}{m}-\sigma}} \|\mathbf{x} - \mathbf{x}^*\| I_{\mathbf{y} \in S_\sigma^c} + \zeta I_{\mathbf{y} \in S_\sigma} \quad (7)$$

where S_σ^c is the complement of S_σ .

Thus, we only need to show inequality (7). For a positive constant c and $\sigma < 1/m$, define

$$S_\sigma = \{\mathbf{y}[t] : \|\mathbf{x}[t] - \mathbf{x}^*\| \leq cK^\sigma\}. \quad (8)$$

Note that $T_i[t] \leq T_{\max}$, and $\|\mathbf{x}[t] - \mathbf{x}^*\| \leq cK^\sigma$ implies that

$$\sum_i x_i[t - s] \leq \sum_i x_i^* + ncK^\sigma + nT_{\max}\hat{\eta}, \quad \text{for all } 0 \leq s \leq T_{\max}.$$

Thus, S_σ is a finite set. Also, it is easy to see that if $\mathbf{y}[t] \in S_\sigma$, there exists $0 < \zeta < \infty$ such that $E[\Delta W_t(\mathbf{y})] < \zeta$. Now, consider $\mathbf{y}[t] \notin S_\sigma$, define the event χ_0^t such that

$$\chi_0^t := \left\{ \max_i \sum_{j=1}^{T_{\max}} a_i[t - j] \leq A \right\}$$

and events χ_l^t for $l = 1, 2, \dots$ such that

$$\chi_l^t := \left\{ \max_i \sum_{j=1}^{T_{\max}} a_i[t - j] = A + l \right\}.$$

Then, we can rewrite $E[\Delta W_t(\mathbf{y})]$ as follows:

$$E[\Delta W_t(\mathbf{y})] = \sum_{l=0}^{\infty} E[\Delta W_t(\mathbf{y}) | \chi_l^t] p(\chi_l^t).$$

For convenience, we also use $\{y\}^M$ to denote $\min\{y, M\}$. Then, along the lines of the proof of [6, Th. 1], it can be shown that there exists $B_d > 0$, which is independent of K and $\mathbf{x}[t]$, such that

$$E[\Delta W_t(\mathbf{y})] \leq \sum_{i=1}^n \Delta x_i[t] \left(\left\{ \frac{\alpha_i K}{(x_i[t - T_i])^m} \right\}^M - \mu_i^* \right) + B_d \quad (9)$$

$$= \sum_{i=1}^n \Delta x_i[t] \left(\left\{ \frac{\alpha_i K}{(x_i[t - T_i])^m} \right\}^M - \left\{ \frac{\alpha_i K}{(x_i[t])^m} \right\}^M \right) \quad (10)$$

$$+ \sum_{i=1}^n \Delta x_i[t] \left(\left\{ \frac{\alpha_i K}{(x_i[t])^m} \right\}^M - \mu_i^* \right) + B_d \quad (11)$$

where $\Delta x_i[t] = x_i[t] - x_i^*$. Define $G(K) := (11)$ and $H(K) := (10)$. To prove inequality (7), we will show the following three facts. The first one is that there exists $\delta_d > 0$ such that for all events χ_l^t

$$G(K) \leq -\frac{\delta_d}{K^{\frac{1}{m}-\sigma}} \|\mathbf{x}[t] - \mathbf{x}^*\|. \quad (12)$$

Second, when χ_0^t happens, there exists $\delta_0 > 0$ such that

$$p(\chi_0^t) |H(K)| \leq p(\chi_0^t) \frac{\delta_0}{K^\xi} \|\mathbf{x}[t] - \mathbf{x}^*\|. \quad (13)$$

The last one is that there exists $\delta_1 > 0$ such that

$$\sum_{l=1}^{\infty} p(\chi_l^t) |H(K)| \leq \frac{\delta_1}{K^\xi} \|\mathbf{x}[t] - \mathbf{x}^*\|. \quad (14)$$

If all three inequalities hold, we have

$$\begin{aligned} E[\Delta W_t(\mathbf{y})] &\leq G(K) + p(\chi_0^t) H(K) + \sum_{l=1}^{\infty} p(\chi_l^t) H(K) \\ &\leq -\left(\frac{\delta_d}{K^{\frac{1}{m}-\sigma}} - \frac{\delta_0 + \delta_1}{K^\xi} \right) \|\mathbf{x}[t] - \mathbf{x}^*\|. \end{aligned}$$

Further, if $\xi > (1/m) - \sigma$, then for $K > ((2\delta_0 + \delta_1)/\delta_d)^{1/(\xi + \sigma - 1/m)}$, we have

$$\frac{\delta_d}{2K^{\frac{1}{m}-\sigma}} - \frac{\delta_0 + \delta_1}{K^\xi} > 0$$

which implies

$$E[\Delta W_t(\mathbf{y})] \leq -\left(\frac{\delta^*}{K^{\frac{1}{m}-\sigma}} \right) \|\mathbf{x}[t] - \mathbf{x}^*\|$$

and the inequality (7) holds with $\delta^* = \delta_d/2$.

Now, we prove (12), (13) and (14). We will first show (12). The proof is similar to the proof of [6, Th. 1]. Here, we consider a general m instead of $m = 1$. Define σ as follows:

$$\sigma = \begin{cases} \lfloor \frac{1}{m} \rfloor, & \text{if } m \leq 1 \text{ and } \frac{1}{m} \text{ is not an integer} \\ \frac{1}{m} - \frac{1}{2}, & \text{if } m \leq 1 \text{ and } \frac{1}{m} \text{ is an integer} \\ \frac{1}{2m}, & \text{if } m > 1. \end{cases}$$

From the previous definition, we have that $0 < (1/m) - \sigma < \min\{\sigma, 1\}$. Note that we have

$$(x_i[t] - x_i^*) \left(\left\{ \frac{\alpha_i K}{(x_i[t])^m} \right\}^M - \mu_i^* \right) \leq 0 \quad \text{for all } i.$$

Letting $i_0 = \arg \max_i |x_i[t] - x_i^*|$, then

$$G(K) \leq -|x_{i_0}[t] - x_{i_0}^*| \left| \left\{ \frac{\alpha_{i_0} K}{(x_{i_0}[t])^m} \right\}^M - \mu_{i_0}^* \right| + B_d.$$

Now, if $\left\{ \frac{\alpha_{i_0} K}{(x_{i_0}[t])^m} \right\}^M = M$, since $M > 2\hat{\eta}$ from the definition of M , we have

$$\left| \left\{ \frac{\alpha_{i_0} K}{(x_{i_0}[t])^m} \right\}^M - \mu_{i_0}^* \right| = M - \mu_{i_0}^* > \hat{\eta}.$$

Otherwise, if $\left\{ \frac{\alpha_{i_0} K}{(x_{i_0}[t])^m} \right\}^M < M$, then

$$\left| \left\{ \frac{\alpha_{i_0} K}{(x_{i_0}[t])^m} \right\}^M - \mu_{i_0}^* \right| = \mu_{i_0}^* \left| \left(\frac{x_{i_0}^*}{x_{i_0}[t]} \right)^m - 1 \right|.$$

Because

$$x_{i_0}[t] = \begin{cases} x_{i_0}^* - |x_{i_0}[t] - x_{i_0}^*| \geq 0, & \text{if } x_{i_0}[t] - x_{i_0}^* \leq 0 \\ x_{i_0}^* + |x_{i_0}[t] - x_{i_0}^*| \geq 0, & \text{if } x_{i_0}[t] - x_{i_0}^* \geq 0 \end{cases}$$

and

$$\begin{aligned} & \left| \left(\frac{x_{i_0}^*}{x_{i_0}^* - |x_{i_0}[t] - x_{i_0}^*|} \right)^m - 1 \right| \\ & \geq \left| \left(\frac{x_{i_0}^*}{x_{i_0}^* + |x_{i_0}[t] - x_{i_0}^*|} \right)^m - 1 \right| \end{aligned}$$

we can obtain that

$$\mu_{i_0}^* \left| \left(\frac{x_{i_0}^*}{x_{i_0}[t]} \right)^m - 1 \right| \geq \mu_{i_0}^* \left| 1 - \frac{1}{(1+\epsilon)^m} \right|$$

where

$$\epsilon = \frac{c\mu_{i_0}^{*1/m}}{\sqrt{n}\alpha_{i_0}^{1/m}} K^{\sigma - \frac{1}{m}} > 0$$

and the inequality holds because $x_{i_0}^* = ((\alpha_{i_0} K)/(\mu_{i_0}^*))^{1/m}$ and $cK^\sigma \leq \|\mathbf{x}[t] - \mathbf{x}^*\| \leq \sqrt{n}|x_{i_0}[t] - x_{i_0}^*|$. Further, since $(1+\epsilon)^m \geq 1+m\epsilon$, we have

$$\mu_{i_0}^* \left| 1 - \frac{1}{(1+\epsilon)^m} \right| \geq \mu_{i_0}^* \left| 1 - \frac{1}{1+m\epsilon} \right| = \frac{\mu_{i_0}^* m\epsilon}{1+m\epsilon}.$$

It is easy to see that for large K , we will have $(\mu_{i_0}^* m\epsilon)/(1+m\epsilon) < \hat{\eta}$. Thus, for large K , we have

$$\begin{aligned} & \left| \left\{ \frac{\alpha_{i_0} K}{(x_{i_0}[t])^m} \right\}^M - \mu_{i_0}^* \right| \geq \frac{\mu_{i_0}^* m\epsilon}{1+m\epsilon} \\ \text{and} \\ & G(K) \leq -|x_{i_0}[t] - x_{i_0}^*| \left(\frac{\mu_{i_0}^* m\epsilon}{1+m\epsilon} - \frac{B_d}{|x_{i_0}[t] - x_{i_0}^*|} \right) \\ & \leq -|x_{i_0}[t] - x_{i_0}^*| \\ & \quad \times \left(\frac{\mu_{i_0}^*}{\left(\frac{\sqrt{n}}{mc} \left(\frac{\alpha_{i_0}}{\mu_{i_0}^*} \right)^{\frac{1}{m}} K^{\frac{1}{m} - \sigma} + 1 \right)} - \frac{B_d}{c\sqrt{n}K^\sigma} \right). \end{aligned}$$

Because $(1/m) - \sigma \leq \sigma$, by choosing sufficiently large c , we can find a positive constant δ and \bar{K} such that for any $K \geq \bar{K}$

$$G(K) \leq -\frac{\hat{\delta}}{K^{\frac{1}{m} - \sigma}} |x_{i_0}[t] - x_{i_0}^*| \leq -\frac{\delta_d}{K^{\frac{1}{m} - \sigma}} \|\mathbf{x}[t] - \mathbf{x}^*\|$$

where $\delta_d = \hat{\delta}/\sqrt{n}$.

Next, we consider (13). It is the case that the arrivals are upper bounded by A . From the assumption, we have $A \geq T_{\max} \hat{\eta}$, so $|x_i[t] - x_i[t - T_{\max}]| \leq A$ and

$$\begin{aligned} & \left| \left\{ \frac{\alpha_i K}{(x_i[t - T_i])^m} \right\}^M - \left\{ \frac{\alpha_i K}{(x_i[t])^m} \right\}^M \right| \\ & \leq M - \frac{\alpha_i K}{\frac{\alpha_i K}{M} \left(1 + \left(\frac{A^m M}{\alpha_i K} \right)^{\frac{1}{m}} \right)^m}. \end{aligned}$$

Since $((A^m M)/(\alpha_i K))^{(1/m)}$ is small when K is large, and $(1+\epsilon)^m \leq 1+2m\epsilon$ for small enough ϵ , we can conclude that for sufficiently large K

$$\begin{aligned} & \left| \left\{ \frac{\alpha_i K}{(x_i[t - T_i])^m} \right\}^M - \left\{ \frac{\alpha_i K}{(x_i[t])^m} \right\}^M \right| \\ & < \frac{2mAM^{\frac{1+m}{m}}}{\alpha_i^{\frac{1}{m}}} \frac{1}{K^{\frac{1}{m}}}. \quad (15) \end{aligned}$$

Let $\alpha_{\min} = \min_i \alpha_i$, then there exists $\delta_0 = 2n\alpha_{\min}^{-(1/m)} mAM^{(1+m/m)}$ and $\xi = (1/m)$ such that

$$|H(K)| \leq \frac{\delta_0}{K^\xi} \|\mathbf{x}[t] - \mathbf{x}^*\|.$$

Finally, we consider the complement of χ_0^t , which is denoted as $\chi_0^{t,c}$, and derive inequality (14). Now the arrivals are not upper bounded and can be arbitrarily large. From assumption (4), the probabilities of these events is very small. So we can still obtain an upper bound for $\sum_{i=1}^n (x_i[t] - x_i^*) | \{ (\alpha_i K)/(x_i[t - T_i]) \}^M - \{ (\alpha_i K)/(x_i[t]) \}^M |$. Now, suppose χ_l^t occurs ($l \geq 1$), similar to (15), we can get

$$\begin{aligned} & \left| \left\{ \frac{\alpha_i K}{x_i[t - T_i]} \right\}^M - \left\{ \frac{\alpha_i K}{x_i[t]} \right\}^M \right| \\ & < 2\alpha_i^{-\frac{1}{m}} mM^{\frac{1+m}{m}} (A+l) \frac{1}{K^\xi} \end{aligned}$$

and

$$\begin{aligned} & \sum_{l=1}^{\infty} p(\chi_l^t) |H(K)| \\ & \leq \|\mathbf{x}[t] - \mathbf{x}^*\| \sum_{l=1}^{\infty} 2n\alpha_{\min}^{-\frac{1}{m}} m M^{\frac{1+m}{m}} \frac{A+l}{K^\xi} p(\chi_l^t). \end{aligned}$$

Under assumption (4), we can further obtain

$$\begin{aligned} & \sum_{l=1}^{\infty} 2n\alpha_{\min}^{-\frac{1}{m}} m M^{\frac{1+m}{m}} \frac{1}{K^\xi} (A+l) p(\chi_l^t) \\ & \leq 2n\alpha_{\min}^{-\frac{1}{m}} m M^{\frac{1+m}{m}} \frac{1}{K^\xi} (h-2) \frac{1}{A^{h-2}} = \frac{\delta_1}{K^\xi} \end{aligned}$$

where $\delta_1 = 2n\alpha_{\min}^{-(1/m)} m M^{(1+m/m)} (h-2)(1/A^{h-2})$.

We have proved that inequalities (12), (13) and (14) hold, and it is easy to see that $\xi > 1/m - \sigma$. Thus, when $K > 4(\delta_0 + \delta_1)^2/\delta_d^2$, inequality (7) holds with $\delta^* = \delta_d/2$, and the lemma follows. ■

Recall that $x_i^* = ((\alpha_i K)/(\mu_i^*))^{(1/m)}$, so from the previous lemma, we have

$$E \left[\frac{|x_i[\infty] - x_i^*|}{x_i^*} \right] \leq E \left[\frac{\|\mathbf{x}[\infty] - \mathbf{x}^*\|}{x_i^*} \right] \leq \frac{\bar{c}(\mu^*)^{1/m}}{\alpha_i^{1/m} K^\sigma}$$

which converges to zero when K goes to infinity. Thus, the mean of $x_i[\infty]$ concentrates around x_i^* for large K , from which we can show that weighted- m fairness can be approximated. This is stated in the next theorem which is the main result of this note.

Theorem 3: Consider the combined scheduling-congestion control algorithm defined by (2)-(6). The steady-state service rate vector $\mu[\infty]$ satisfies the following: For any $\epsilon > 0$

$$\lim_{K \rightarrow \infty} P(|\mu[\infty] - \mu^*| \geq \epsilon) = 0.$$

Proof: Using the Markov inequality, Lemma 2 yields for any $\epsilon > 0$

$$P \left(\frac{1}{K^{1/m}} |x_i[\infty] - x_i^*| > \epsilon \right) \leq \frac{\bar{c}}{\epsilon K^\sigma}.$$

Further, since

$$\mu[t] \in \arg \max_{\mu \in C_s[t]} \sum_{i=1}^n x_i[t] \eta_i = \arg \max_{\mu \in C_s[t]} \sum_{i=1}^n \frac{x_i[t]}{K^{1/m}} \eta_i$$

we can conclude

$$\lim_{K \rightarrow \infty} P(|\mu[\infty] - \mu^*| \geq \epsilon) = 0$$

and the network is weighted m -fair according to Theorem 1. ■

Theorem 3 allows us to conclude that even in the presence of delays, the network will approach the weighted- m fairness. Note that, from inequality (7), when K is large, $\{\mathbf{y}[t]\}$ is positive recurrent and the system is stable. Actually if we are only concerned with the stability of the system, inequality (7) is much stronger than what is necessary to prove the stability. In fact, we can show that for any $K > 0$, the Markov chain is positive recurrent. Define the $S_{\bar{X}}$:

$$S_{\bar{X}} = \left\{ \mathbf{y}[t] : \sum_i x_i[t] \leq \bar{X} \right\}. \quad (16)$$

Clearly, $S_{\bar{X}}$ is a finite set. Stability of the system is established by following theorem.

Theorem 4: For any $K < 0$, there exists positive numbers ζ , \bar{X} and δ such that

$$E[\Delta W_t(\mathbf{y})] \leq -\delta \sum_{i=1}^n x_i[t] I_{\mathbf{y} \in S_{\bar{X}}} + \zeta I_{\mathbf{y} \in S_{\bar{X}}}$$

where $S_{\bar{X}}$ is defined as (16). Hence, the Markov chain $\{\mathbf{y}[t]\}$ is positive recurrent.

Proof: We omit the proof here because of lack of the space. Please refer to [20, Th. 4] for the proof of $m = 1$. The case of general weighted- m fairness is similar. ■

From Theorem 3, we see that fairness can only be achieved when $K \rightarrow \infty$. However, Theorem 4 assures that we are guaranteed at least stability for all K .

IV. CONCLUSION

In this note, we have shown that the algorithm (5) and (6) is stable even in the presence of heterogeneous delays and when K is large, the network will approach the weighted- m fairness. When delays are not negligible in some situations, our result reinforces the result that the combination of queue-length-based scheduling and congestion control is a good distributed fair resource allocation scheme.

REFERENCES

- [1] J. Aikat, J. Kaur, F. D. Smith, and K. Jeffay, "Variability of TCP round-trip times," in *Proc. ACM Sigcomm, Internet Measure. Conf.*, Oct. 2003, pp. 279–284.
- [2] M. Andrews, K. Kumaran, K. Ramanan, A. Stolyar, R. Vijayakumar, and P. Whiting, "Scheduling in a queueing system with asynchronously varying service rates," Bell Laboratories, Tech. Rep., 2000.
- [3] M. Armony and N. Bambos, "Queueing dynamics and maximal throughput scheduling in switched processing systems," Stanford Univ., Stanford, CA, Tech. Rep. Netlab-2001-09/01, 2001.
- [4] S. Asmussen, *Applied Probability and Queues*. New York: Springer-Verlag, 2003.
- [5] R. Buche and H. J. Kushner, "Control of mobile communication systems with time-varying channels via stability methods," *IEEE Trans. Autom. Control*, vol. 49, no. 11, pp. 1954–1962, Nov. 2004.
- [6] A. Eryilmaz and R. Srikant, "Fair resource allocation in wireless networks using queue-length-based scheduling and congestion control," in *Proc. IEEE Infocom*, Miami, FL, Mar. 2005, pp. 1794–1803.
- [7] A. Eryilmaz, R. Srikant, and J. Perkins, "Stable scheduling policies for fading wireless channels," in *Proc. IEEE Int. Symp. Information Theory*, Apr. 2005, vol. 13, pp. 411–424.
- [8] F. P. Kelly, A. Maulloo, and D. Tan, "Rate control in communication networks: Shadow prices, proportional fairness and stability," *J. Oper. Res. Soc.*, vol. 49, pp. 237–252, 1998.
- [9] R. Leelaahakriengkrai and R. Agrawal, "Scheduling in multimedia wireless networks," in *Proc. ITC*, Salvador de Bahia, Brazil, 2001.
- [10] X. Lin and N. Shroff, "The impact of imperfect scheduling on cross-layer rate control in wireless networks," in *Proc. IEEE Infocom*, Miami, FL, Mar. 2005, pp. 1804–1814.
- [11] S. H. Low and D. E. Lapsley, "Optimization flow control I: Basic algorithm and convergence," *IEEE/ACM Trans. Networking*, vol. 7, no. 6, pp. 861–875, Dec. 1999.
- [12] J. Mo and J. Walrand, "Fair end-to-end window-based congestion control," *IEEE/ACM Trans. Networking*, vol. 8, no. 5, pp. 556–567, Oct. 2000.
- [13] M. Neely, E. Modiano, and C. Li, "Fairness optimal stochastic control for heterogeneous networks," in *Proc. IEEE Infocom*, Miami, FL, Mar. 2005, pp. 1723–1734.
- [14] M. J. Neely, E. Modiano, and C. E. Rohrs, "Dynamic power allocation and routing for time varying wireless networks," in *Proc. IEEE Infocom*, Apr. 2003, pp. 745–755.
- [15] S. Shakkottai and A. Stolyar, , Yu. M. Suhov, Ed., "Scheduling for multiple flows sharing a time-varying channel: The exponential rule," in *Analytic Methods in Applied Probability*, ser. AMS Trans. 2, 207. Providence, RI: AMS, 2002.

- [16] R. Srikant, *The Mathematics of Internet Congestion Control*. Boston, MA: Birkhäuser, 2004.
- [17] A. Stolyar, "Maximizing queueing network utility subject to stability: Greedy primal-dual algorithm," *Queueing Syst.*, vol. 50, Aug. 2005.
- [18] L. Tassiulas and A. Ephremides, "Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks," *IEEE Trans. Autom. Control*, vol. 37, no. 12, pp. 1936–1948, Dec. 1992.
- [19] H. Yaiche, R. R. Mazumdar, and C. Rosenberg, "A game-theoretic framework for bandwidth allocation and pricing in broadband networks," *IEEE/ACM Trans. Networking*, vol. 8, no. 5, pp. 667–678, Oct. 2005.
- [20] L. Ying, R. Srikant, A. Eryilmaz, and G. Dullerud, "Distributed fair resource allocation in cellular networks in the presence of heterogeneous delays," in *Proc. WiOpt*, 2005, pp. 96–105.

Relay-Based Identification of a Class of Nonminimum Phase SISO Processes

Somanath Majhi

Abstract—A set of general expressions is derived from a single symmetrical relay feedback test for identification of a class of process transfer functions. Using these expressions the parameters of open loop stable nonminimum phase transfer function models may be obtained from simple measurements made on the limit cycle. For comparison, the conventional describing function based identification formulae are presented. Fourier series based curve fitting with the options of nonlinear least squares method and trust-region algorithm is used to measure limit cycle parameters in the presence of measurement noise. Examples are given to illustrate the value of the proposed method.

Index Terms—Identification, measurement noise, nonminimum phase, relay feedback.

I. INTRODUCTION

Study and analysis of relay feedback system is a classical field. Relays in electromechanical systems and simple models of dry friction motivated extensive analysis of relay control systems earlier. A vast collection of analysis methods and applications of relay control systems can be found in [1]–[4]. More recent applications include delta-sigma modulations and automatic tuning of proportional–integral–derivative (PID) controllers. As for the automatic tuning of PID controllers implemented in many industries, the idea is to determine some points on the Nyquist curve of a stable open-loop plant by measuring the frequency of oscillation induced by a relay feedback [5]. However, to achieve better performance, a number of controller tuning methods use model based designs such as [6]–[9]. Therefore, to apply these tuning rules a proper parametric model of the process is needed. This has resulted in several studies [10]–[15] that aim to identifying suitable models using a relay feedback test. However, it is observed that while many of the relay autotuning methods use the describing function (DF) approximation to a relay in their analysis, a very few have applied the exact relay feedback expressions but involving more complexity. Despite the

Manuscript received July 6, 2005; revised March 1, 2006. Recommended by Associate Editor A. Garulli.

The author is with the Department of Electronics and Communication Engineering, Indian Institute of Technology Guwahati, North Guwahati, Assam PIN-781039, India (e-mail: smajhi@iitg.ernet.in).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TAC.2006.886520

apparent success in industrial applications, the basic relay autotuning method can be improved significantly to obtain a more accurate estimate of the plant parameters. The reason is that the estimated ultimate gain and ultimate frequency derived from the describing function are only an approximation for information at the critical frequency.

Recently, Panda and Xu [16] have derived exact expressions for relay feedback responses for estimation of the transfer functions of a class of processes. However, their method needs to locate exactly the starting point on the relay response that lies between the two zero crossing points. Under realistic condition when some measurement noise is present it is very difficult to identify the starting point. Further, it is not clear how one can choose correct model structures for unknown parameter estimation for their type 3 class of nonmonotonic responses since the responses used to have same starting points. Therefore, their method may not work properly under nonmonotonic relay responses or presence of measurement noise. Measurement noise is a common problem in almost all process industries. The accuracy of the relay experiment is adversely effected if there is measurement noise present in the system output, in particular spurious switching of the relay can be experienced. In this note, the exact analysis of the relay limit cycle is used and a set of simple analytical expressions are derived recursively. The expressions can be used with measurements taken from a single symmetrical relay test for identification of a class of nonminimum phase processes. The proposed method for model identification is tested against measurement noise.

II. PROCESS MODEL

Typical nonminimum phase transfer function models in process control are often assumed to be stable time delay process models having no zeros and possessing a single or multiple poles or stable process models having a right half plane zero and multiple poles. Therefore, the following transfer function model having some generality:

$$G(s) = \frac{K(-T_0s + 1)e^{-\theta s}}{(T_1s + 1)^p} = \frac{K(-\lambda)^p(-T_0s + 1)e^{-\theta s}}{(s - \lambda)^p} \quad (1)$$

is considered for this identification problem where $p = 1, 2, 3 \dots$, and $\lambda = -1/T_1$. The two typical nonminimum phase process models can be obtained by setting $\{T_0 = 0, \theta \neq 0\}$ and $\{\theta = 0, T_0 \neq 0\}$. When $\theta = T_0 = 0$, the process models represent a minimum phase system. The reason why the transfer function models are chosen will be apparent in the development that follows. For each model the number of unknown parameters is three and each relay test supplies three data.

When the p th order process model is expressed in the Jordan canonical form, its state equation constants \mathbf{A} , \mathbf{b} , and \mathbf{c} become

$$\mathbf{A} = \begin{bmatrix} \lambda & 1 & 0 & \cdot & \cdot & \cdot & \cdot & 0 \\ 0 & \lambda & 1 & 0 & \cdot & \cdot & \cdot & 0 \\ 0 & 0 & \lambda & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & 0 & \cdot & \cdot & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & 0 & \lambda & 1 \\ 0 & 0 & 0 & \cdot & \cdot & 0 & 0 & \lambda \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \\ 1 \end{bmatrix} \quad (2)$$

$$\mathbf{c} = K(-\lambda)^p [(1 - T_0\lambda) \quad -T_0 \quad 0 \quad \dots \quad 0].$$

III. THE RELAY FEEDBACK METHOD

Consider the relay to be an ideal one with output amplitude $\pm h$. The state and output equations for the process plant are

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{b}u(t - \theta) \quad (3)$$

$$y(t) = \mathbf{c}\mathbf{x}(t) \quad (4)$$