



HAL
open science

Fairness in MIMD Congestion Control Algorithms

Eitan Altman, Konstantin Avrachenkov, Balakrishna Prabhu

► **To cite this version:**

Eitan Altman, Konstantin Avrachenkov, Balakrishna Prabhu. Fairness in MIMD Congestion Control Algorithms. [Research Report] RR-5312, INRIA. 2004, pp.29. inria-00070688

HAL Id: inria-00070688

<https://inria.hal.science/inria-00070688>

Submitted on 19 May 2006

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

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

Fairness in MIMD Congestion Control Algorithms

E. Altman — K. E. Avrachenkov — B. J. Prabhu

N° 5312

Septembre 2004

Thème COM

 ***rapport
de recherche***

Fairness in MIMD Congestion Control Algorithms

E. Altman , K. E. Avrachenkov , B. J. Prabhu

Thème COM — Systèmes communicants
Projets Maestro

Rapport de recherche n° 5312 — Septembre 2004 — 29 pages

Abstract: We study fairness among sessions sharing a common bottleneck link, where one or more sessions use a multiplicative increase multiplicative decrease (MIMD) algorithm. Losses or congestion signals occur when the capacity is reached but could also be initiated before that. Both synchronized as well as non-synchronized losses are considered. In the non-synchronized case, only one session suffers a loss at a time. Two models are then considered to determine which source loses a packet: a *rate dependent* model in which the loss probability of a session is proportional to its rate at the congestion instant, and the *independent loss rate* model. We first study how two MIMD sessions share the capacity in the presence of general combinations of synchronized and non-synchronized losses. We show that, in the presence of rate dependent losses, the capacity is fairly shared whereas rate independent losses provides high unfairness. We then study inter protocol fairness: how the capacity is shared in the presence of synchronized losses among sessions some of which use additive increase multiplicative decrease (AIMD) protocols whereas the others use MIMD protocols.

Key-words: Fairness, MIMD, synchronized and non-synchronized losses, stochastic stability

Sur l'équité dans les algorithmes de contrôle de congestion A.M.D.M

Résumé : Nous étudions l'équité lorsque plusieurs sources utilisant l'algorithme des accroissements multiplicatifs et de la décroissance multiplicative (A.M.D.M.) de TCP partagent un goulot d'étranglement. On dit que les pertes de paquets sont synchronisées lorsque toutes les sources subissent une perte au même instant. Chiu et Jain ont démontré que A.M.D.M était inéquitable lorsque les pertes de paquets étaient synchronisées. Dans un premier temps, nous étudions le partage de bande passante entre deux sources A.M.D.M. lorsque les pertes de paquets sont non-synchronisées. Nous démontrons que la bande passante est partagée équitablement lorsque les pertes de paquets sont dépendantes du débit de la source alors que le partage est inéquitable lorsque les pertes de paquets sont indépendantes du débit de la source. Ensuite, nous étudions l'équité entre les sources A.A.D.M. (des accroissements additifs et de la décroissance multiplicative) et les sources A.M.D.M qui partagent un goulot d'étranglement lorsque le processus de pertes est synchronisé.

Mots-clés : Equité, A.M.D.M, pertes synchronisées et non-synchronisées, stabilité stochastique

1 Introduction

In the Internet, data transfer protocols use different congestion control algorithms to achieve rate control. Until now, the AIMD algorithm was found to provide satisfactory performance. However, in high speed networks, MIMD algorithm (e.g., [1, 2]) has been proposed in order to efficiently utilise the network capacity. Therefore, in the future, situations may arise where different sessions using these two algorithms would compete for the same network resource. The share of the capacity obtained by each of these sessions will depend on the various parameters specific to the algorithms. The sharing of a resource gives rise to the question on how fairly is this resource shared. Fairness issues have been addressed in several previous works. In [3], the authors consider a class of rate control algorithms under the assumption of synchronized control signals, and show that the AIMD algorithm converges to fairness. In [4], the authors consider MIMD algorithms under the more realistic assumption of state dependent losses, and argue that MIMD algorithm also converges to fairness. In [5] the convergence to fairness of the different flavours of TCP are studied both analytically and using simulations. Loguinov *et al.* [6] study the monotonic convergence to fairness for algorithms in rate-based TCP-friendly applications. In [7], the authors mention that for sessions with different round trip times (RTT), Scalable TCP (which uses MIMD algorithm) is extremely unfair. They propose a new algorithm to improve fairness.

Losses or congestion signals (these terms will be used interchangeably) occur when the capacity is achieved but could also be initiated before that. Losses are said to be synchronized when all the sessions sharing the link suffer a loss at the same time instant. In the non-synchronized case, only one session suffers a loss. Both synchronized losses as well as non-synchronized losses are considered as well as their combination. Two types of models are then considered to determine which of the sources loses a packet: (1) the *rate dependent loss model* in which the loss probability of a session is proportional to its rate at the congestion instant [8], and (2) the *independent loss rate model*, in which the session to which a signal is sent is independent of the session's rate.

We first study how two MIMD sessions (with either the same RTT or with different RTTs) share the capacity in the presence of general combinations of synchronized and non-synchronized losses. We show that, in the presence of rate dependent losses, the capacity is fairly shared between the two sessions whereas rate independent losses result in high unfairness even when sessions are symmetric. In the second part we study how the capacity is shared among several sessions, each of which either uses an additive increase multiplicative decrease (AIMD) or a MIMD algorithm in presence of synchronized losses only. We show that the AIMD session obtains a share which is independent of the link capacity, and that the rest of the capacity is utilized by the MIMD session.

The rest of the paper is organised as follows. In Section 2 we present a brief overview of the model and mention the contribution of this work. In the first part of the paper (Sections 3 and 4), we analyze the fairness of two MIMD sessions sharing a common link in the presence of rate dependent losses. In the second part of the paper (Section 5 and Section 6), we study fairness of AIMD and MIMD sessions sharing a common link. We present the conclusions in Section 7.

2 Overview

We use the following notation. A function from \mathbb{R} to \mathbb{R} will be denoted using sans serif font, such as $x(\cdot)$. For example, $x(t)$ denotes a function defined for all real values of t . A function from \mathbb{Z} to \mathbb{R} will be denoted using italic fonts as $x(\cdot)$. Usually, $x(n)$ would be the value of $x(t)$ at the n^{th} sampling instant. A vector a will denote a row vector. Its transpose will be denoted by a' . Also, we use the term session to mean an instance of a given algorithm. The term user will be used for someone who makes use of one or more instances (i.e., sessions) of the same algorithm.

The following model is mainly based on the model in [3]. Consider two flows which share a link of capacity C . Let $x(t) \equiv (x_1(t), x_2(t))$ be the rate vector at time t , where $x_1(t)$ and $x_2(t)$ denote the instantaneous rates of session 1 and of session 2, respectively. The set of feasible rate vectors, $\{(x_1, x_2) | x_1 + x_2 \leq C; x_1, x_2 \geq 0\}$ is shown in Fig. 1. The line

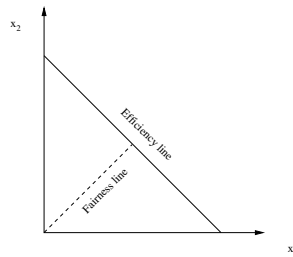


Figure 1: The rate allocation vector

$x_1(t) + x_2(t) = C$ is called the efficiency line. On this line the available capacity is fully utilized. The line $x_1(t) = x_2(t)$ is called the fairness line. On this line both the sessions obtain the same rates, and hence the bandwidth sharing is said to be fair. The sessions react to control signals by adapting their rates in the following way.

$$x(n+1) = \begin{cases} b_I x(n) + a_I & \text{for an increase signal,} \\ b_D x(n) + a_D & \text{for a decrease signal,} \end{cases}$$

where a_I , b_I , a_D and b_D are constants and the sampling is done just after the control instants. (In TCP Reno for example, arrivals of ACKs can be considered as increase signals and arrivals of duplicated ACKs or a time-out is considered to be a decrease signal.) In [3] it has been argued that for convergence to the fairness line, the increase algorithm has to be multiplicative and additive (i.e., $a_I > 0$ and $b_I \geq 1$) and the decrease algorithm has to be multiplicative ($b_D < 1$ and $a_D = 0$). The authors assume that the control signals are synchronized for both the sessions, and that the control signal is the same for both the sessions. That is, both the sessions receive control signals at the same instant and both

of them either increase, or decrease their rates simultaneously¹. Indeed, under these two assumptions, the rate vector for MIMD algorithm stays on a line joining the origin to the initial rate vector, and hence does not converge to fairness. In [4], the authors show that under more realistic assumption of rate dependent control signal, MIMD algorithm also converges to the fairness line. In the first part of the paper, we first show that, for sessions with the same RTT, MIMD algorithm converges to fairness when the control signals are rate dependent. We obtain the expressions for the long term fairness index, the rate of convergence to the steady state distribution and the mean time to achieve fairness. In [7] it was argued that for sessions with different RTTs, Scalable TCP (or, MIMD algorithm) is extremely unfair. We show that, even for sessions with different RTTs, a certain degree of fairness can be achieved by introducing sufficient number of asynchronous losses. We then show, through simulations, that the results obtained for two sessions also hold for n sessions.

In the second part of this paper, we consider several sessions sharing a common link on which losses are due to buffer overflow and are, therefore, synchronized. In [5] such a scenario was considered for session only using AIMD algorithm. Loguinov *et al.* [6] provide fairness and packet-loss scalability analysis and simulations for session using more general binomial algorithms. In [7], fairness issues were considered for high-speed networks. The RTT-fairness was compared for different proposals of TCP in high-speed networks. Here, too, the fairness was studied between sessions using the same algorithm. However, we consider a heterogeneous scenario where different sessions may use different algorithms. This type of scenarios may be of interest in the future when, for example, sessions using Scalable TCP and standard TCP will share the same link. We analyze the equilibrium throughput and the window size of the sessions and compare them with simulations. We note that, as was pointed out in [6], the window-based notation can be converted to a rate-based notation using the relation $x(n) = w(n) \frac{MTU}{RTT}$, where $w(n)$ is the window at the n^{th} sampling instant, M is the packet size in bits and RTT is the round trip time in seconds. Therefore, we shall use the rate-based notation in the first part of the paper, and the window-based notation in the second part of the paper.

3 Fairness in MIMD sessions (equal RTTs)

We consider two sessions which share a link of capacity C . At time t , the rates obtained by the two sessions are denoted by $\mathbf{x}(t) \equiv (x_1(t), x_2(t))$. At each control instant, the controller sends a control signal to each source. This signal either informs on no congestion (a 0 signal) or sends a congestion (a 1) signal. In the absence of congestion, the sources increase their rate exponentially, i.e.,

$$x_i(t + \tau) = \alpha^{\tau/\tau_0} \cdot x_i(t), \quad i = 1, 2,$$

where τ_0 is the time constant (for example, the RTT) for the sessions, and $\alpha > 1$ is the increase factor. The above formulation is a continuous time equivalent of a multiplicative

¹These assumptions are validated in [9] by simulations for AIMD versions of TCP, with approximately the same RTTs, and we validate them later for MIMD and AIMD versions of TCP (see, e.g., Fig. 5 and its discussion).

Table 1: Reaction to control signals

control vector	$x_1(t_j+)$	$x_2(t_j+)$
(0, 0)	$x_1(t_j)$	$x_2(t_j)$
(0, 1)	$x_1(t_j)$	$\beta \cdot x_2(t_j)$
(1, 0)	$\beta \cdot x_1(t_j)$	$x_2(t_j)$
(1, 1)	$\beta \cdot x_1(t_j)$	$\beta \cdot x_2(t_j)$

algorithm in which, for every RTT without congestion signal, the sender multiplies the window by a factor of α . This can be seen by substituting $t = n\tau_0$. The control signals to the two sources are assumed to be synchronised, i.e., they receive control signals at the same instant. However, unlike the model in [3], the two sources can receive different control signals. Let $\beta < 1$ be the decrease factor. Let the j^{th} control signal be received at time t_j . Then, the four possibilities for the rate vector, $x(t_j+)$, just after t_j , are given in Table 1. The source continues with the increase algorithm on the reception of 0 signal. On the other hand, when a source receives a 1 signal, it instantaneously reduces its rate. We assume that whenever the link capacity is attained then either a synchronized or a non-synchronized loss occurs. Furthermore, non-synchronized losses may occur before attaining the capacity. We note that if there were only synchronized losses then $x_i(n) = x_i(0)$, $i = 1, 2$, for all integers n so that any initial unfair sharing would remain forever. Therefore, if MIMD protocols are used then it is essential to provide a stream of non synchronized congestion signals using, for example, some queue management scheme.

3.1 Instantaneous throughput ratio process

We now study the instantaneous throughput ratio process. It is shown to be a Markov chain with a countable state space; we shall show that this chain could be stable or unstable depending on the asynchronous loss process.

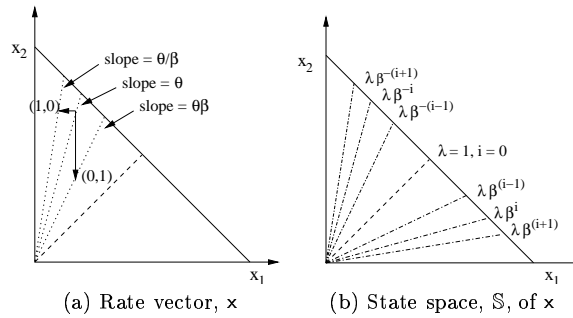


Figure 2: Geometric interpretation.

In Figure 2(a), we show the geometric interpretation of the response to the different control signals. Let θ be the slope of the line joining the origin and the current vector, $\mathbf{x}(t)$. That is,

$$\theta = \frac{x_2(t)}{x_1(t)}.$$

If there are no control signals in the interval $(t, t + \tau)$ then the rate vector at time $(t + \tau)$ will be $\alpha^{\tau/\tau_0}(x_1(t), x_2(t))$. This vector also lies on the line with slope θ . If a $(1, 1)$ signal was generated at t , then the vector after the response, $\mathbf{x}_1(t+)$, is $\beta(x_1(t-), x_2(t-))$. This vector also lies on the line with slope θ . Therefore, the rate vector remains on the line of its slope as long as the control vector is either $(0, 0)$ or $(1, 1)$, or there are no control signals. However, it can be seen that the rate vector moves to the line with slope θ/β when a control signal of $(1, 0)$ is generated. Similarly, the rate vector moves to the line with slope $\theta\beta$ when a control vector of $(0, 1)$ is generated. Therefore, given an initial slope of θ_0 , the slope of the line along which the rate vector lies just after the n^{th} control signal can be written as $\theta_n = \theta_0\beta^i$ for some $i \in \mathbb{Z}$, where \mathbb{Z} is the set of all integers. For any given initial slope, θ_0 , we can find a unique $\lambda \in (\beta^{1/2}, \beta^{-1/2})$ and $j \in \mathbb{Z}$ such that θ_0 can be expressed in terms of λ as $\theta_0 = \lambda\beta^j$. For convenience, we will define θ_n in terms of λ . The state space of θ_n is a countably infinite state space defined by

$$\mathbb{S} = \{\lambda\beta^i, \forall i \in \mathbb{Z}\}.$$

A geometric interpretation of \mathbb{S} is shown in Fig. 2(b). We note that the line $\lambda = 1, i = 0$ is the fairness line. The continuous time increase and instantaneous decrease of the algorithm allows us to obtain the above formulation.

In the rest of this section, we assume that $\lambda = 1$. This assumption is equivalent to saying that the initial vector has a slope of β^i . If λ is not equal to 1 then, by any combination of control signals, the instantaneous rate vector can only get close to the fairness line by getting to λ . The rate vector can not, however, be on the fairness line.

Let $s(t) = \{i : i \in \mathbb{Z}\}$ be the process which denotes that the rate vector at time t lies on a line with slope β^i . We embed this process at instants of arrival of the control signals. Let t_n denote the time instant when the n^{th} congestion occurs. Then, $s_n = s(t_n+) = s(n)$.

We shall assume below that congestion signals occur only when capacity is reached, and that at these events a non-synchronized loss is generated with a probability $0 < \epsilon < 1$. We shall later show that the qualitative results obtained carry also to the case of congestion signals sent also before capacity is attained.

3.2 Rate dependent loss model

Next, we consider the rate dependent loss model in which the loss probability of a session is proportional to the session's rate at the congestion instant. In particular, we assume that the probability of loss for session i at the n^{th} loss instant is given by $\frac{x_i(n)}{x_1(n)+x_2(n)}$.

Proposition 3.1 $\{s_n, n \geq 0\}$, is a discrete state-space Markov chain with transition probabilities given by

$$\begin{aligned}
 s_{n+1}|(s_n = i) &= \begin{cases} i+1 & w.p. \frac{\epsilon \beta^{|i|}}{1+\beta^{|i|}} \\ i-1 & w.p. \frac{\epsilon}{1+\beta^{|i|}}, \\ i & w.p. 1-\epsilon \end{cases}, \quad i > 0 \\
 s_{n+1}|(s_n = 0) &= \begin{cases} +1 & w.p. \epsilon/2 \\ -1 & w.p. \epsilon/2 \\ 0 & w.p. 1-\epsilon \end{cases} \\
 s_{n+1}|(s_n = i) &= \begin{cases} i-1 & w.p. \frac{\epsilon \beta^{|i|}}{1+\beta^{|i|}} \\ i+1 & w.p. \frac{\epsilon}{1+\beta^{|i|}}, \\ i & w.p. 1-\epsilon \end{cases}, \quad i < 0
 \end{aligned} \tag{1}$$

Remark 3.1 So far we have excluded losses that occur before capacity is reached. Consider now the case that non-synchronized losses are generated occasionally independently of the rates of the connections. Synchronized losses occur whenever capacity is reached. Let t_n be the n^{th} instant in which a non-synchronized loss occurs. Then Proposition 3.1 still holds with $\epsilon = 1$. Therefore the stability conditions we shall obtain will also hold for this scenario.

We are interested in finding the steady state distribution of the Markov chain, and the mean first passage time to the state $i = 0$ starting from a random state. The state $i = 0$ corresponds to the fairness line. Therefore, the mean first passage time from a random state to the state $i = 0$ gives an indication of the mean time before the rate vector reaches the fairness line. In general, the first passage time to state $i = 0$ gives the first passage time to λ . If $\lambda \neq 1$, s_n cannot be on the fairness line, and so the above performance measure corresponds to first passage time to the state closest to the fairness line for this particular process.

The fairness index at the n^{th} control instant, F_n , is defined as follows

$$F_n = \frac{1}{2} \frac{(x_1(n) + x_2(n))^2}{x_1(n)^2 + x_2(n)^2}. \tag{2}$$

We can write F_n in terms of β^i as follows

$$F_n = \frac{1}{2} \sum_{i=-\infty}^{\infty} P(s_n = i) \frac{(1 + \beta^i)^2}{1 + \beta^{2i}}.$$

The long term fairness index, F_∞ , can be obtained by taking the limit $n \rightarrow \infty$. We assume that the process s_n converges to its stationary limit, s_∞ .² Then, F_∞ can also be expressed

²We establish this later. In cases that the Markov chain is null recurrent (Subsection 3.3) we shall understand the steady state distribution to correspond to the compactification of the state space in which $-\infty$ and ∞ are added to the state space and each will have a stationary probability of 0.5.

as

$$F_\infty = \frac{1}{2} + \sum_{i=-\infty}^{\infty} P(s_\infty = i) \frac{\beta^i}{1 + \beta^{2i}}. \quad (3)$$

Stability. The existence of the limiting distribution, s_∞ , can be ensured by proving that the Markov chain s is positive recurrent. A Markov chain is positive recurrent if it satisfies the Foster's criterion [10] which is stated below.

Theorem 3.1 (Foster) *An irreducible Markov chain s , on a countable state \mathbb{Z} , is ergodic if and only if there exists a positive function $f(\alpha)$, $\alpha \in \mathbb{Z}$, a number $\mu > 0$ and a finite set A such that*

$$E[f(s_{n+1}) - f(s_n) | s_n = i] \leq -\mu, \quad i \notin A, \quad (4a)$$

$$E[f(s_{n+1}) | s_n = i] < \infty, \quad i \in A. \quad (4b)$$

Let $f(i) = |i|$, $i \in \mathbb{Z}$ and $A = \{0\}$. Let Δf_i be define as $\Delta f_i := E[f(s_{n+1}) - f(s_n) | s_n = i]$. First, we show that condition (4a) is satisfied. For $i \neq 0$, from (2), we have

$$\begin{aligned} \Delta f_i &= \epsilon \frac{\beta^{|i|}}{1 + \beta^{|i|}} (|i| + 1) \\ &\quad + \epsilon \left(1 - \frac{\beta^{|i|}}{1 + \beta^{|i|}} \right) (|i| - 1) + (1 - \epsilon) |i| - |i| \\ &= \epsilon \left(2 \frac{\beta^{|i|}}{1 + \beta^{|i|}} - 1 \right) \leq -\epsilon \left(1 - 2 \frac{\beta}{1 + \beta} \right). \end{aligned}$$

which is strictly negative for any $\beta \in [0, 1)$. Therefore, condition (4a) is satisfied. To check for (4b), for $i = 0$ we have

$$E[f(s_{n+1}) | s_n = i] = \frac{\epsilon}{2} |-1| + \frac{\epsilon}{2} |1| = \epsilon < \infty.$$

Therefore, the Markov chain s satisfies the conditions of Theorem 3.1 and, hence, is positive recurrent.

3.3 The independent loss rate model

If the session to which a congestion signal is sent is independent of the rates then the transition probabilities become

$$s_{n+1} | (s_n = i) = \begin{cases} i - 1 & \text{w.p. } \frac{\epsilon}{2} \\ i + 1 & \text{w.p. } \frac{\epsilon}{2}, \\ i & \text{w.p. } 1 - \epsilon \end{cases} \quad \forall i \quad (5)$$

We show that this results in instability. We use the following theorem from [10].

Theorem 3.2 *For an irreducible Markov chain s to be null recurrent, it suffices that there exist two functions $f(x)$ and $\psi(x)$, $x \in \mathbb{Z}$, and a finite subset $A \in \mathbb{Z}$, such that the following conditions hold:*

- 1) $f(x) \geq 0$, $\psi(x) \geq 0$, $\forall x \in \mathbb{Z}$.
- 2) For some positive α, γ , with $1 < \alpha \leq 2$,

$$f(x) \leq \gamma[\psi(x)]^\alpha, \forall x \in \mathbb{Z}.$$

- 3) $\lim_{x_i \rightarrow \infty} \psi(x_i) = \infty$ and $\sup_{x \notin A} f(x) > \sup_{x \in A} f(x)$. 4)

- a) $E[f(s_{n+1}) - f(s_n) | s_n = x] \geq 0$, $x \notin A$;
- b) $E[\psi(s_{n+1}) - \psi(s_n) | s_n = x] \leq 0$, $x \notin A$;
- c) $\sup_{x \in X} E[|\psi(s_{n+1}) - \psi(s_n)|^\alpha | s_n = x] = C, < \infty$.

Let $f(x) = x^2$, $\psi(x) = |x|$, $A = \{0\}$, $\gamma = 1$ and $\alpha = 2$. Conditions 1, 2 and 3 of the above theorem are satisfied with these assumptions. For condition 4(a),

$$\Delta f_i = \frac{\epsilon}{2}((x+1)^2 - x^2) + \frac{\epsilon}{2}((x-1)^2 - x^2) = \epsilon > 0.$$

For condition 4(b),

$$\Delta \psi_i = \frac{\epsilon}{2}((|x|+1) - |x|) + \frac{\epsilon}{2}((|x|-1) - |x|) = 0.$$

For condition 4(c),

$$\begin{aligned} E[|\psi(s_{n+1}) - \psi(s_n)|^2 | s_n = x] = \\ \frac{\epsilon}{2}((|x|+1) - |x|)^2 + \frac{\epsilon}{2}((|x|-1) - |x|)^2 = \epsilon < \infty. \end{aligned}$$

Hence, the Markov chain s is null recurrent. Since s is null recurrent, and the discrete state space has two accumulation points, $(C, 0)$ and $(0, C)$, on the line $x_1 + x_2 = C$, the probability of being in any small vicinity of each point is $\frac{1}{2}$. The mean time to go from one extreme to another will be ∞ and, therefore, one connection will get the whole capacity. This suggests that rate independent losses are not sufficient to improve the fairness whereas rate dependent losses can indeed provide a fair share of the capacity.

In the sequel, we thus focus on the rate dependent loss model.

3.4 Steady state distribution

The Markov chain s is positive recurrent, and, therefore, the steady state distribution $P(s_\infty = i)$ exists. Since s_n is symmetric about the state 0, we can consider a Markov

chain, $\{y_n, n \geq 0\}$, on the state space $\{0, 1, 2, \dots\}$, in order to obtain the steady state distribution of s . The transition probabilities at the n^{th} control instant for this random walk are given by

$$y_{n+1}|(y_n = i) = \begin{cases} i + 1 & \text{w.p. } \epsilon q_i \\ i - 1 & \text{w.p. } \epsilon(1 - q_i) \\ i & \text{w.p. } 1 - \epsilon \end{cases}, i > 0,$$

$$y_{n+1}|(y_n = i) = \begin{cases} i + 1 & \text{w.p. } \epsilon \\ i & \text{w.p. } 1 - \epsilon \end{cases}, i = 0,$$

where $q_i = \frac{\beta^i}{1 + \beta^i}$. Let y_∞ denote the steady state process to which y_n converges. Let p_i denote the probability of y_∞ being in state i . Then

Proposition 3.2

$$p_i = p_0 \prod_{j=0}^{i-1} \frac{q_j}{1 - q_{j+1}}. \quad (6)$$

p_0 can be obtained from the equation $\sum_{i=0}^{\infty} p_i = 1$.

Proof 3.1 Let $p_n = 0, \forall n < 0$. The balance equation for this walk can be written as

$$\begin{aligned} p_i &= (1 - \epsilon)p_i + \epsilon q_{i-1}p_{i-1} + \epsilon(1 - q_{i+1})p_{i+1} \\ &= q_{i-1}p_{i-1} + (1 - q_{i+1})p_{i+1} \end{aligned} \quad (7)$$

These are the balance equations for a birth-death process with state dependent transition probabilities. The solution for this type of process is known to be of the form in (6)[11].

Note that p_i can also be written as $p_i = p_0 \beta^{i(i-1)/2} (1 + \beta^i)$. Since $\beta^i \rightarrow 0$ as $i \rightarrow \infty$, the tail of p_i can be seen to decrease as β^{i^2} which is a very fast decrease. In particular, this means that the process is around the fairness line most of the time.

The steady state distribution of s_∞ can be obtained from the following relations

$$P(s_\infty = i) = \begin{cases} \frac{1}{2}P(y_\infty = |i|), & \text{for } i \neq 0 \\ P(y_\infty = |i|), & \text{for } i = 0. \end{cases}$$

The long term fairness index, F_∞ , can be computed numerically using Eqns. (3) and (6). In Table 2 we give the fairness index for different value of β .

Table 2: Fairness index for different values of β .

β	0.95	0.875	0.75	0.6	0.5	0.1
F_∞	0.987	0.97	0.942	0.91	0.88	0.777

3.5 Convergence to steady state distribution

The second largest eigenvalue of a matrix gives the rate of convergence to the steady state distribution. Therefore, we can get an indication of the rate of convergence of the Markov chain y by looking at the eigenvalues of its transition probability matrix, P . P can be written as

$$P = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & \dots \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ \vdots \end{matrix} & \begin{pmatrix} 1 - \epsilon & \epsilon & 0 & \dots \\ \epsilon(1 - q_1) & 1 - \epsilon & \epsilon q_1 & \dots \\ 0 & \epsilon(1 - q_2) & 1 - \epsilon & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \end{matrix} \quad (8)$$

Let ζ_i denote the i^{th} eigenvalue of P such that $\zeta_i \geq \zeta_j$ for $i < j$, and $\zeta = 1$. We can obtain the following lower bound on ζ_i .

Proposition 3.3 For $i > 0$, $1 - \epsilon \leq \zeta_i < 1$.

Proof 3.2 Since y_n is irreducible, the multiplicity of eigenvalues at 1 is 1. Therefore, $\zeta_i < 1$ for $i > 0$.

We can rewrite P as

$$P = (1 - \epsilon)I + \epsilon A \quad (9)$$

where I is the identity matrix and A is a transition matrix of a pure birth-death process with up transition probability q_i and down transition probability $1 - q_i$. A is a stochastic matrix and, therefore, all its eigenvalues belong to the interval $[0, 1]$. Let μ_i be the i^{th} eigenvalue of A and let v_i be the corresponding left eigenvector. Then, from (9) we get

$$v_i P = (1 - \epsilon)v_i + \epsilon v_i A = ((1 - \epsilon) + \epsilon \mu_i)v_i.$$

Therefore, v_i is also the left eigenvector of P , and the corresponding eigenvalue is $(1 - \epsilon) + \epsilon \mu_i$. Since $\mu_i \geq 0$, we get the inequality $\zeta_i \geq 1 - \epsilon$.

Therefore, $1 - \epsilon$ gives a lower bound on the rate of convergence of the Markov chain to the steady state.

3.6 Mean first passage time

In this section we compute the mean first passage time to the state 0 starting from a random state. This gives us an estimate of the first time the rate vector reaches the fairness line starting from a given initial random state. We note that the Markov chain s is a birth-death process which is symmetric about the state 0. If the initial state is positive, the Markov chain will stay in the set of positive states before visiting state 0. Similarly, if the initial state is negative, the Markov chain will stay in the set of negative states before visiting state 0. Therefore, we can obtain the mean first passage time to state 0 for s by obtaining the mean first passage time to state 0 for y .

Let $p = (p_0 \ p_1 \ \dots)$ be the steady state probability vector of P , the transition probability matrix of y , as given by (8). Its i^{th} component, p_i , is given by Eqn. (6). Let $m = (m_1 \ m_2 \ \dots)$ denote the mean first passage time vector with m_j , $j \geq 1$ denoting the mean first passage time to the state 0 starting from state j .

Proposition 3.4 $m_i, i \geq 2$ can be obtained from the following recursion: $m_0 = 0$,

$$m_1 = \frac{1}{\epsilon} \frac{1 - p_0}{p_0}, \quad (10)$$

$$m_{i+1} = \frac{\epsilon m_i - \epsilon \frac{1}{1+\beta^i} m_{i-1} - 1}{\epsilon \frac{\beta^i}{1+\beta^i}}. \quad (11)$$

Proof 3.3 Let P_1 be the transition probability matrix conditioned on y not being in state 0. We can rewrite P as

$$P = \begin{pmatrix} 0 & \dots \\ 1 - \epsilon & a \\ \vdots & b & P_1 \end{pmatrix}, \quad (12)$$

where the vector a is given as $a = (\epsilon \ 0 \ \dots)$ and the vector b is defined as $b = (1 - \beta \ 0 \ \dots)'$.

The vector m satisfies the equation [12]

$$(I - P_1)m' = \underline{1}', \quad (13)$$

where $\underline{1}$ is vector of ones. Since P_1 is a tridiagonal matrix, we can rewrite the above equation as

$$m_{i+1} = \frac{\epsilon m_i - \epsilon \frac{1}{1+\beta^i} m_{i-1} - 1}{\epsilon \frac{\beta^i}{1+\beta^i}}, \quad i \geq 2, \quad (14)$$

with the definition $m_0 = 0$.

Let the vector r be defined as $r = (p_1 \ p_2 \ \dots)$. Since p is the steady state probability vector of P , we have $pP = p$. We can rewrite the above equation using (12).

$$(p_0 \ r) \begin{pmatrix} 1 - \epsilon & a \\ b & P_1 \end{pmatrix} = (p_0 \ r). \quad (15)$$

Solving for r , we get $r(I - P_1) = p_0 \cdot a$. Multiplying the above equation by m' , and substituting for $(I - P_1) \cdot m'$ from Eqn.(13), we obtain

$$r \cdot \underline{1}' = p_0 \cdot a \cdot m'.$$

Substituting $r \cdot \underline{1}' = 1 - p_0$ in the above equation, we obtain

$$\frac{1 - p_0}{p_0} = a \cdot m'.$$

The vector a has ϵ in its first column and 0 elsewhere. Therefore, the above equation reduces to (10). (14) together with (10) give the desired recursion.

We note that the steady state probabilities are independent of ϵ whereas the mean first passage times are inversely proportional to ϵ .

4 Fairness in MIMD sessions (unequal RTTs)

In this section we assume that the two sessions have different time constants. Let τ_1 and τ_2 be the time constants of session 1 and of session 2, respectively. The rate evolution for session i in the absence of control signals can be written as

$$x_i(t + \tau) = x_i(t)\alpha^{\tau/\tau_i}, i = 1, 2.$$

We now make the following transformation

$$z(t) = \log \left[\frac{x_2(t)}{x_1(t)} \right]. \quad (16)$$

We consider again both synchronized losses as well as non-synchronized ones. In the absence of non-synchronized losses, the evolution of $z(t)$ becomes

$$z(t + \tau) = z(t) + \gamma\tau,$$

where $\gamma = \log[\alpha] \left(\frac{1}{\tau_2} - \frac{1}{\tau_1} \right)$. If a control signal arrives at t , then $z(t+)$ can be written as

$$z(t+) = \begin{cases} z(t) & \text{signal is either } (0,0) \text{ or } (1,1) \\ z(t) + \log[\beta] & \text{signal is } (0,1) \\ z(t) - \log[\beta] & \text{signal is } (1,0) \end{cases}.$$

The evolution in time of z is shown in Figure 3. Here we assume that γ is positive, i.e.,

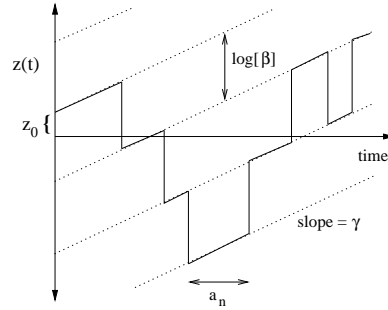


Figure 3: Evolution in time of $z(t)$. $\tau_2 < \tau_1$.

$\tau_2 < \tau_1$. Therefore, if there were only synchronized signals then there would be a drift in time towards $+\infty$. This suggests that the rate for session 1 would approach 0. Similarly, if

γ were negative, i.e., $\tau_2 > \tau_1$, there would be a drift towards $-\infty$ suggesting that the rate for session 2 would approach 0. It can also be seen that if the slope were to be 0, i.e., τ_1 were to be equal to τ_2 , the dotted lines would be parallel to the time axis. In this case $z(t)$ would remain constant in the absence of non-synchronized losses which was also observed in the previous section. We conclude that some buffer management scheme that creates non-synchronized losses is necessary in order to have some fairness.

We therefore assume in the sequel that some buffer management scheme occasionally creates rate dependent non-synchronized losses. Let $\{a_n, n \geq 0\}$ denote the time between the n^{th} and the $(n+1)^{\text{th}}$ such signals. The process a_n is assumed to be i.i.d. Let $z_n = z(t_n)$ be the process embedded just before the arrival of a non-synchronized control signal. The probability that the n^{th} such signal is for session i can be rewritten as $\frac{1}{e^{z_n} + 1}$ for session 1, and $\frac{e^{z_n}}{e^{z_n} + 1}$ for session 2.

Proposition 4.1 *The process $\{z_n, n \geq 0\}$ is a Markov chain with state space \mathbb{R} , and it follows the recursive equation*

$$z_{n+1} = z_n + \gamma a_n + c_n, \quad (17)$$

where c_n is defined as

$$c_n = \begin{cases} -\log[\beta] & \text{w.p. } \frac{1}{e^{z_n} + 1} \\ +\log[\beta] & \text{w.p. } \frac{e^{z_n}}{e^{z_n} + 1} \end{cases}.$$

4.1 Stability

We are interested in knowing the arrival rates for which the process $z(t)$ does not have a drift towards ∞ as $t \rightarrow \infty$. Let $u_n = \gamma a_n$, and let $U(\cdot)$ and v be the distribution function and the mean, respectively, of u_n . We assume that the density function, $\frac{dU(\cdot)}{du}$, of u is a non-increasing function. Let $b = -\log[\beta]$. We note that b is a positive number since β is less than 1.

Proposition 4.2 *The Markov chain z defined by (17) is positive recurrent if*

$$v < b. \quad (18)$$

Proof 4.1 *To show the positive recurrence of the Markov chain we use the following theorem from [13].*

Theorem 4.1 ([13]) *For some “small”³ set $W \in \mathbb{B}(\mathbb{R})$, some constant $h < \infty$, $\mu > 0$, and an extended real-valued function $V : \mathbb{R} \rightarrow [0, \infty]$, z is stable if*

$$\Delta V(x) := \int_{\mathbb{R}} P(x, dy) V(y) - V(x) \leq -\mu + h \mathbf{1}_W(x), \quad (19)$$

³ A set W is called a “small” set if there exists a measure ϕ , $\phi(\mathbb{R}) > 0$, such that

$$P(x, A) \geq \phi(A), \quad x \in C, A \in \mathbb{B}(\mathbb{R}).$$

where $P(x, \cdot)$ is the one step transition probability matrix of z .

To check for the drift condition of this theorem, we consider $V(y) = |y|$, and the set $W = [-b, b]$. The LHS of (19) becomes

$$\begin{aligned}
\Delta V(x) + |x| &= \int_{y \in \mathbb{R}} \frac{1}{1 + e^x} |y| dU(y - (x + b)) \\
&+ \int_{y \in \mathbb{R}} \frac{e^x}{1 + e^x} |y| dU(y - (x - b)) \\
&= \int_{y=(x+b)}^{\infty} \frac{1}{1 + e^x} |y| dU(y - (x + b)) \\
&+ \int_{y=(x-b)}^{\infty} \frac{e^x}{1 + e^x} |y| dU(y - (x - b)) \\
&= \int_0^{\infty} \frac{1}{1 + e^x} |y + (x + b)| dU(y) \\
&+ \int_0^{\infty} \frac{e^x}{1 + e^x} |y + (x - b)| dU(y) \\
&\leq \int_0^{\infty} \frac{1}{1 + e^x} (|y| + |x + b|) dU(y) \\
&+ \int_0^{\infty} \frac{e^x}{1 + e^x} (|y| + |x - b|) dU(y) \\
\Delta V(x) + |x| &\leq v + \int_0^{\infty} \frac{1}{1 + e^x} |x + b| dU(y) \\
&+ \int_0^{\infty} \frac{e^x}{1 + e^x} |x - b| dU(y). \tag{20}
\end{aligned}$$

For $x \in W$, (20) can be rewritten as

$$\begin{aligned}
\Delta V(x) + |x| &\leq v + \int_0^{\infty} \frac{1}{1 + e^x} (|x| + |b|) dU(y) \\
&+ \int_0^{\infty} \frac{e^x}{1 + e^x} (|x| + |b|) dU(y) \\
\Delta V(x) &\leq v + b < \infty.
\end{aligned}$$

For $x \in W^c$, (20) can be rewritten as

$$\begin{aligned}
\Delta V(x) + |x| &\leq v + \int_0^{\infty} \frac{1}{1 + e^{|x|}} (|x| + b) dU(y) \\
&+ \int_0^{\infty} \frac{e^{|x|}}{1 + e^{|x|}} (|x| - b) dU(y) \\
\Delta V(x) &\leq v + \frac{1 - e^{|x|}}{1 + e^{|x|}} b.
\end{aligned}$$

Let x^* be the value of x for which

$$\Delta V(x) \leq v + \frac{1 - e^x}{1 + e^x} b = -\mu.$$

Then, $x^* = \log \frac{b+v+\mu}{b-(v+\mu)}$. For $b > v + \mu$, there exists a $x^* \in \mathbb{R}$ for which the $\Delta V(x) \leq -\mu$. If x^* is less than b then the drift condition is satisfied for the $W = [-b, b]$. However, if x^* is greater than b then we can consider the set $W = [-b, x^*]$. Hence, for $W = [-b, \max(b, x^*)]$, the drift condition (19) is satisfied. It follows from Lemma 4.1 that W is indeed a small set.

Lemma 4.1 For any d such that $-b \leq d < \infty$, the set $W = [-b, d]$ is a “small” set.

Proof 4.2 For $x \in W$,

$$\begin{aligned} P(x, A) &= \frac{1}{1 + e^x} \int_{y \in A} dU(y - (x + b)) \\ &+ \frac{e^x}{1 + e^x} \int_{y \in A} dU(y - (x - b)) \\ &\geq \frac{e^x}{1 + e^x} \int_{y \in A} dU(y + b - x) \end{aligned} \quad (21)$$

Since $\frac{dU(u)}{du}$ was assumed to be a non increasing function in u . Then, $\int_{y \in A} dU(y + b - x)$ is non decreasing function in x . Also, $\frac{e^x}{1 + e^x}$ is an increasing function in x . Since $x \geq -b$, we can rewrite (21) as

$$P(x, A) \geq \frac{e^{-b}}{1 + e^{-b}} \int_{y \in A} dU(y + 2b) \geq \phi(A), \quad (22)$$

where $\phi(A) := \frac{e^{-b}}{1 + e^{-b}} \int_{y \in A} dU(y + 2b)$. Since there exists a measure ϕ such that $P(x, A) \geq \phi(A)$, $x \in W, A \in \mathbb{B}(\mathbb{R})$, a closed and bounded set $W = [-b, d]$ is a small set.

Since W is a “small set” when $b > v$, from Theorem 4.1 we can conclude that, for $b > v$, the Markov chain z is positive recurrent.

We note that $b > v$ is a sufficient condition for positive recurrence. Let $\frac{1}{\lambda}$ be the mean time between losses. Then $\frac{1}{\lambda} = \frac{v}{\gamma}$. From Prop. 4.2, for z to be positive recurrent

$$\lambda > \frac{\gamma}{b}. \quad (23)$$

Therefore, to achieve some fairness, the arrival rate of the losses process has to be greater than $\frac{\gamma}{b}$.

4.2 Simulation results

In this subsection we present the results of simulations. Our objective is to verify the analytical result obtained in Prop. 4.1 which noted that sufficient number of asynchronous losses are required so that sessions with different RTTs can share the capacity fairly. In the simulation scenario, nine Scalable TCP sessions shared a link of 200Mbps. Sessions 1,2 and 3 had a RTT of 50ms. Sessions 4,5 and 6 had a RTT of 90ms, and sessions 7,8 and 9 had a RTT of 140ms. The simulations were performed using *ns-2*(version 2.26)[14]. In Fig. 4, the window size is plotted as a function of time for different values of ϵ (i.e., packet drop probability). We note that $\epsilon = 0$ corresponds to only congestion losses which are seen to be not always synchronous. Therefore, in Fig. 4(a) there are asynchronous as well as synchronous losses even though $\epsilon = 0$. However, during periods of synchronous losses (which have been pointed out in the figure) there is short-term unfairness. Even though there are asynchronous losses due to congestion, the window sizes of the sessions with larger RTTs go to 0. We now induce further asynchronous losses by dropping each packet with probability $\epsilon \neq 0$. In Figs. 4(b) and 4(c) there is a marked improvement in the throughput obtained by sessions with larger RTTs as the loss probability is increased. For small loss probability, there is still some unfairness between sessions with different RTTs. However, for sessions with RTT of 50ms, there is no short-term unfairness as was observed when there were no induced asynchronous losses. For a larger loss probability (i.e., $\epsilon = 0.0003$), sessions share the capacity fairly. This confirms the analytical result which stated that the fairness in MIMD sessions with different RTTs can be achieved by introducing sufficient asynchronous losses. Let η_1 , η_2 and η_3 be the total throughput of sessions with RTT of 50ms, 90ms

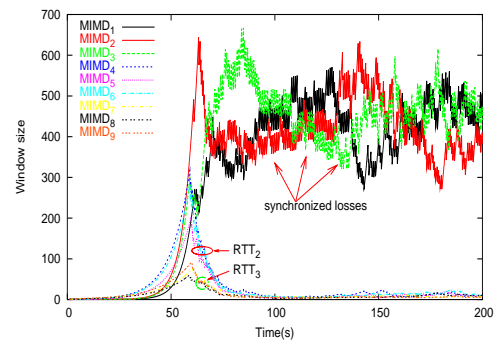
Table 3: Throughput for each RTT class and overall efficiency

ϵ	η_1 (Mbps)	η_2 (Mbps)	η_3 (Mbps)	$\frac{\eta_1 + \eta_2 + \eta_3}{C}$
0	178	2.8	1	0.91
0.00015	148	25.5	7.14	0.905
0.0003	101	48	28.4	0.89

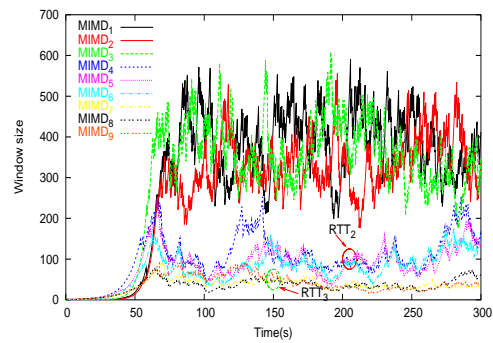
and 140ms, respectively. In Table 3, the values of throughput and the overall efficiency are given. We note that, for $\epsilon = 0.0003$, the ratio of the throughputs of two different classes, $\frac{\eta_i}{\eta_j}$, are almost in proportion of the respective RTTs. Therefore, we can say that a certain degree of fairness has been achieved at the cost of marginal decrease in efficiency.

5 Inter protocol fairness (Same RTT)

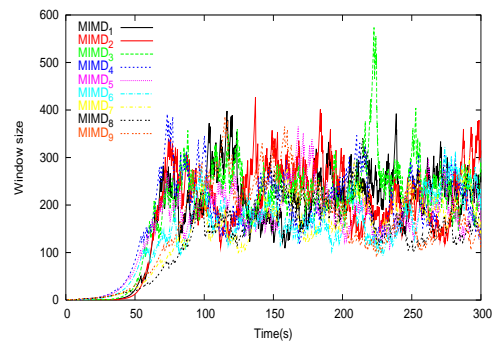
In the second part of this paper, we study the fairness issue when sessions using two different congestion control algorithms share a common link. Recently, Scalable TCP, which uses MIMD algorithm, has been proposed as an enhancement for TCP in high-speed networks. Situations may, therefore, arise in which a user with Scalable TCP shares a link with a user



(a) $\epsilon = 0$.



(b) $\epsilon = 0.00015$.



(c) $\epsilon = 0.0003$.

Figure 4: Window evolution MIMD sessions. $RTT_1 = 50ms$. $RTT_2 = 90ms$. $RTT_3 = 140ms$.

with standard TCP. Specifically, we study the equilibrium behaviour of the window size, and the throughput obtained by a session of each algorithm at equilibrium in the presence of synchronized losses only. We also look at conditions under which a user of one algorithm can obtain a better throughput than a user of the other algorithm. Previous work (e.g., [5], [7]) mainly studied the behaviour of sessions using the same protocol.

In this section, we assume that each session has the same RTT, τ . As mentioned in Section 2, window-based notation is equivalent to rate-based notation. In the rest of this paper, we use the window-based notation since we are interested in obtaining the equilibrium window sizes for the sessions.

5.1 System Model

Consider l sessions which share a link of capacity C bits/s. Each session transmits data using packets of size M bits. Let Λ be the bandwidth-delay product (BDP) network. We assume that the RTT is mainly determined by the propagation delay and, hence, can be considered to be a constant.

Let $\mathbf{x}(t) = (x_1(t) \ x_2(t) \ \dots \ x_l(t))$ denote the vector of window sizes of the l sessions at time t . A synchronized loss (i.e., a loss for each session) is assumed to occur at time t if

$$\sum_{i=1}^l x_i(t) > \Lambda. \quad (24)$$

The above condition is equivalent to saying that a synchronized loss occurs when the total number of outstanding packets in the network exceeds the total number of packets that the network can handle.

Without loss of generality, let sessions $1, 2, \dots, k$ use the MIMD congestion control algorithm and the rest of the $l - k$ sessions use the AIMD congestion control algorithm. In the absence of losses, the two algorithms increase the window in the following way

$$x_i(t + \Delta) = \begin{cases} x(t) \alpha_m^{\Delta/\tau} & \text{for } 1 \leq i \leq k \\ x(t) + \alpha_a \frac{\Delta}{\tau} & \text{for } k + 1 \leq i \leq l, \end{cases} \quad (25)$$

where α_m and α_a are the increase parameters of the MIMD and the AIMD algorithm, respectively. For example, $\alpha_m = 1.01$ for Scalable TCP, and $\alpha_a = 1$ for standard TCP. Let t_n denote the time instant when the n^{th} congestion signal is received. We note that a congestion signal is generated when a synchronized loss occurs. In response to a congestion signal the two algorithms decrease the window in the following way

$$x_i(t_n^+) = \begin{cases} \beta_m x(t_n) & \text{for } 1 \leq i \leq k, \\ \beta_a x(t_n) & \text{for } k + 1 \leq i \leq l, \end{cases}$$

where β_m and β_a are the decrease parameters of the MIMD and the AIMD algorithm, respectively. For example, $\beta_m = 0.875$ for Scalable TCP, and $\beta_a = 0.5$ for standard TCP.

Let $x(n)$ denote the window-size vector embedded just after the n^{th} congestion signal is received, i.e, $x(n) = x(t_n+)$. Let δ_n denote the time between two congestion signals. Since all the sessions are assumed to receive congestion signals at the same instant, we can write the following recursive equation for $x(n)$.

$$x_i(n+1) = \begin{cases} \beta_m x_i(n) \alpha_m^{\delta_n/\tau} & \text{for } 1 \leq i \leq k, \\ \beta_a (x_i(n) + \alpha_a \frac{\delta_n}{\tau}) & \text{for } k+1 \leq i \leq l. \end{cases} \quad (26)$$

5.2 Bandwidth Sharing

The transient behaviour of the window sizes can be obtained by solving Eqns (26) and (24). Given the initial window vector $x(0)$, the time to the first loss t_1 and, hence, $x(1)$ can be computed. This way we can recursively compute $x(n)$. This allows us to obtain the behaviour of the window-size vector and the loss instants before the equilibrium is reached. At equilibrium, δ_n and $x(n)$ will converge to their steady state values denoted by δ^* and ψ , respectively. We are interested in finding the window size, ψ_i , of each session at equilibrium. Then, ψ_i together with δ^* will allow us to obtain the throughput for session i . At equilibrium $x(n)$ would be identical to $x(n+1)$, $x(n+2)$, and so on. Therefore, for each session i , we can obtain ψ_i from Eqn (26) as follows.

$$\psi_i = \begin{cases} \beta_m \psi_i \alpha_m^{\delta^*/\tau} & \text{for } 1 \leq i \leq k, \\ \beta_a (\psi_i + \alpha_a \delta^*/\tau) & \text{for } k+1 \leq i \leq l. \end{cases} \quad (27)$$

The l equations in (27) are fixed point solutions of the corresponding equations in (26). We now have $n+1$ variables in $\psi_i, 1 \leq i \leq n$ and δ^* , and n equations. The final equation can be obtained by noting that a loss occurs when condition (24) is satisfied. Therefore, the $(n+1)^{\text{th}}$ equation is obtained as

$$\sum_{i=1}^k \frac{\psi_i}{\beta_m} + \sum_{i=k+1}^l \frac{\psi_i}{\beta_a} = \Lambda. \quad (28)$$

We note that from any one of the first k equations in (27) we can obtain the value of δ^* . The variables $\psi_i, 1 \leq i \leq k$ cannot be uniquely determined from these k equations. They will depend on the window vector just after the first synchronized loss. This result is equivalent to the result obtained in [3] where the rate vector of symmetric MIMD sessions was dependent on the initial rate vector. However, the equilibrium window size for the AIMD sessions can be uniquely determined from (27). Therefore,

$$\delta^* = \tau \frac{\log[1/\beta_m]}{\log \alpha_m}, \quad (29)$$

$$\psi_i = \alpha_a \frac{\beta_a}{1 - \beta_a} \frac{\log[1/\beta_m]}{\log \alpha_m} \quad k+1 \leq i \leq l. \quad (30)$$

From (28) we can, however, obtain the sum of the equilibrium window sizes of the MIMD sessions.

$$\sum_{i=1}^k \psi_i = \beta_m \Lambda - \frac{\beta_m}{\beta_a} \sum_{j=k+1}^l \psi_j, \quad (31)$$

In order to compute the throughput, η_i , for session i , we divide the time interval τ in slots of length δ^* . We note that, just after a loss instant, the window size for session i is ψ_i . In between two loss instants, the window size for each session increases using the algorithm given by Eqn. (25). Also, in every RTT (i.e., in every slot), session i transfers packets equivalent to its present window size. Therefore, in between two loss instants, the total number of packets that are transferred by session i can be obtained by summing the window sizes during the δ^*/τ RTTs. As before, we can obtain the throughput η_i for each AIMD session whereas we can obtain the total throughput, η_m , for all the MIMD sessions.

$$\eta_m = \frac{M}{\delta^*} \sum_{i=1}^k \psi_i \frac{\alpha_m^{\lfloor \frac{\delta^*}{\tau} \rfloor + 1} - 1}{\alpha_m - 1} \quad (32)$$

$$\eta_i = \frac{M\alpha_a}{\tau} \left(\psi_i + \alpha_a \frac{\lfloor \frac{\delta^*}{\tau} \rfloor + 1}{2} \right) \quad k+1 \leq i \leq l. \quad (33)$$

We note that the throughput expressions are approximate since the number of packets transferred in an RTT is an integer whereas ψ_i can take non integer values. Also, the number of packets transferred in the RTT in which a loss occurs may not be equal to $\lfloor \psi_i \rfloor$.

We can make the following observations from Eqns. (29)- (33). The equilibrium value of the time between two loss instants, δ^* , is independent of the parameters of the AIMD algorithm. It is determined by the RTT, τ , and the parameters of the MIMD algorithm only. The equilibrium window size of the AIMD sessions depends only on the increase and decrease parameters of the two algorithms. Also, the AIMD sessions have the same equilibrium window behaviour and, hence, obtain the same throughput. The rest of the capacity is utilized by the MIMD sessions.

Simulations: We now compare these observations with simulations performed using *ns-2* (version 2.26). Unless stated otherwise, the simulation had the same set of parameters. The MIMD sessions used Scalable TCP, and the AIMD sessions used TCP New Reno. The packet size, M , for each session was set to 1040 bytes (1000 bytes of data + 40 bytes of header). The propagation delay, σ , was taken to be 100ms. The increase and decrease parameters for the two algorithms were set to $\alpha_m = 1.01$, $\alpha_a = 1.0$, $\beta_m = 0.75$, and $\beta_a = 0.5$. Since the ψ_i for AIMD increases with decrease in β_m , we set β_m to a value smaller than its recommended value so that the AIMD sessions also obtain a certain throughput. From Figs. 5(a) (3 MIMD sessions and 3 AIMD sessions) and 5(b) (6 MIMD sessions and 6 AIMD sessions), we note that the AIMD sessions indeed converge to the same equilibrium window size whereas the equilibrium window size of an MIMD session depends on its window just before the first synchronized loss. The ψ_i for AIMD sessions remains the same even though the link capacity is increased from 200Mbps to 300Mbps and the total number of sessions is increased

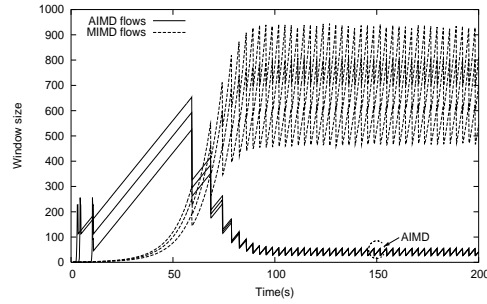
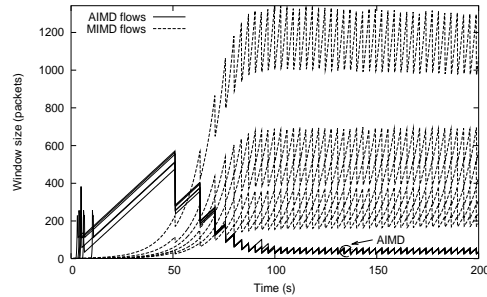
(a) $C = 200\text{Mbps}$. $l = 6$. $k = 3$.(b) $C = 300\text{Mbps}$. $l = 12$. $k = 6$.

Figure 5: Window evolution of sessions.

Table 4: Several MIMD and several AIMD sessions.

	Link Speed (Mbps)	δ (s)	η_m (Mbps)	η_a (Mbps)	ψ_a (packets)
$l = 6$. $k = 3$.	200	2.83 (3)	164(151)	3.5(3.4)	29 (27)
$l = 12$. $k = 6$.	300	2.83 (3)	238.6(218)	3.5(3.4)	29 (27)

from six to twelve. Let η_a and ψ_a denote the throughput and the equilibrium window size of any one of the AIMD sessions, respectively. In Table 4, the analytical and simulation values of δ^* , η_m , η_a , and ψ_a are given. The simulation values are given in parentheses. As predicted in the analysis, the equilibrium window size and the throughput of the AIMD sessions remains unchanged even when the capacity is increased from 200Mbps to 300Mbps, and the total number of sessions is increased from six to twelve.

5.3 Throughput Comparison

We now study the scenario where one MIMD user and AIMD user share the same link. We note that each user can initiate several sessions of the same algorithm. We are interested in knowing the conditions under which the AIMD user can obtain better throughput than the MIMD user.

First, we consider the case in which each user initiates only one session. In such a scenario, the window size and the throughput of each session is obtained from Eqns. (30)-(32) with $l = 2$ and $k = 1$.

From (32) and (33), as $\Lambda \rightarrow \infty$ (i.e., $C \rightarrow \infty$), the ratio of the throughputs, $\frac{\eta_2}{\eta_1}$, goes to 0. This suggests that in high-speed networks, the MIMD user will get most of the capacity. On the other hand, if the BDP of the network is small, the MIMD user will obtain lower throughput compared to the AIMD session.

Proposition 5.1 *Let Λ_l denote a threshold BDP below which an AIMD session will get better throughput compared to a MIMD session. The threshold value, Λ_l , is given by*

$$\Lambda_l = \psi_2 \left(\frac{\alpha_a \delta^*}{2\tau} \left(\frac{\delta^*}{\tau} + 1 \right) \kappa + \frac{1}{\beta_a} \right), \quad (34)$$

where $\kappa = \frac{\alpha_m - 1}{\alpha_m - \beta_m}$.

Proof 5.1 *The above relation can be obtained using the fact $\eta_1 \leq \eta_2$ together with Eqn. (31).*

The value of Λ_l depends only on the increase and decrease parameters of the two algorithms. Table 5 gives the values of Λ_l for different β_m with $\alpha_m = 1.01$ and $\alpha_a = 1$.

Table 5: Λ_l for different values of β_m .

β_m	0.875	0.75	0.5
Λ_l	47.34	106.6	282.73

In Fig. 6(a), the window evolution is plotted for the two sessions for $C = 13$ Mbps and $\beta_m = 0.5$. The BDP, Λ , is less than the Λ_l . The AIMD algorithm obtains a better throughput in this case. In the next set of simulations, we set β_m to its recommended value of 0.875. In Fig. 6(b), the corresponding window evolution is plotted. The effect of increasing β_m is to reduce the share of the AIMD session. In Table 6, a comparison of the values obtained from analysis and simulations is presented. A good match is observed between the analysis and simulations.

From (33), it was observed that the throughput obtained by each AIMD session remains constant whereas the total throughput of the MIMD sessions increases with increase in capacity. An AIMD user may want to obtain throughputs similar to a MIMD user. In this case, the AIMD user may open several sessions in order to improve its observed throughput.

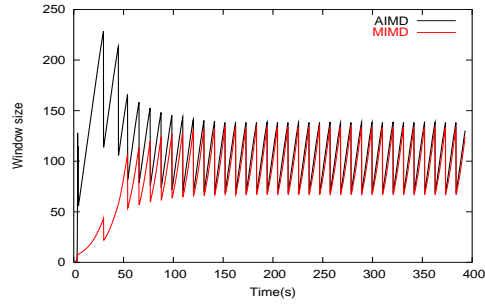
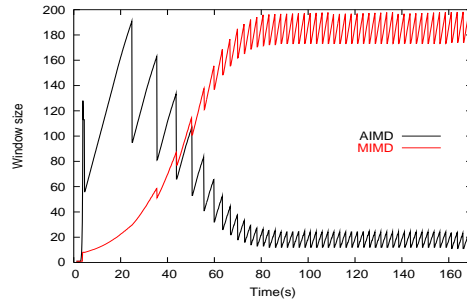
(a) $C = 13\text{Mbps}$. $\beta_m = 0.5$. $\tau = 140\text{ms}$.(b) $C = 10\text{Mbps}$. $\beta_m = 0.875$. $\tau = 140\text{ms}$.

Figure 6: Window evolution for one MIMD session and one AIMD session.

Table 6: One MIMD and one AIMD session. $\tau = 140\text{ms}$.

	Link Speed (Mbps)	η_1 (Mbps)	η_2 (Mbps)	ψ_1 (packets)	ψ_2 (packets)
$\beta = 0.5$	3	2.22 (2.36)	0.6 (0.52)	66.2 (70.7)	13.41 (11)
	5	4.05 (4.1)	0.75 (0.71)	96.82 (98)	13.41 (12)
	10	8.86 (8.74)	0.92 (0.86)	173.39 (173)	13.41 (12)
$\beta = 0.875$	13	4.73 (5.19)	4.92 (5.64)	69.09 (67)	69.6 (69)
	15	6.08 (6.62)	5.04 (5.65)	86.59 (84)	69.6 (68)
	30	16.64 (16.95)	5.47 (5.86)	217.83 (211)	69.6 (69)

Since each AIMD session gets the same throughput independent of the number of AIMD sessions (assuming there is sufficient capacity), an AIMD user can improve its observed throughput by opening multiple sessions.

Proposition 5.2 *The smallest number of sessions, ν , with which an AIMD user will obtain a better global throughput compared to single MIMD user is given by*

$$\nu = \left\lceil \frac{\beta_m}{h\kappa} \left(\Lambda - \frac{\psi_a}{\beta_a} \right) \right\rceil, \quad (35)$$

where $h = \psi_a \alpha_a \sum_{j=0}^{\lfloor \frac{\delta^*}{\tau} \rfloor} j$, ψ_a is the equilibrium window of any one of the AIMD sessions and is defined in (30), and κ is as defined in (34).

Note that ν depends only on Λ and the increase and decrease parameters of the two algorithms. Table 7 gives the value of ν for different values of Λ for $\beta_m = 0.875$.

Table 7: ν for different values of Λ .

$\beta_m =$	Λ	100	500	1000	10000	50000
0.5	ν	3	18	37	372	1863
$\beta_m =$	Λ	100	500	1000	10000	50000
0.875	ν	1	3	7	71	358

Similar to the AIMD user, a MIMD user may also try to improve its observed throughput by opening several sessions. Since, from (33), the AIMD user will get a throughput independent of the number of MIMD sessions, the observed throughput of an MIMD user will not improve by opening several sessions. This result is in contrast to the result obtained in (35) where we noted that an AIMD user can improve its observed throughput by opening several sessions.

6 Inter protocol fairness (Different RTTs)

In this section we study the effect of having a different RTT for each session on the equilibrium window behaviour. The notation used and the scenario is the same as in Sec. 5. We assume that there exists a BDP, Λ , such that there is a synchronized loss when condition (24) is satisfied. Let τ_i be the RTT of session i . Then, we can rewrite (27) as follows.

$$\psi_i = \beta_m \psi_i \alpha_m^{\delta^*/\tau_i}, \quad 1 \leq i \leq k, \quad (36)$$

$$\psi_i = \beta_a (\psi_i + \alpha_a \delta^*/\tau_i), \quad k+1 \leq i \leq l. \quad (37)$$

The expressions for throughput are as follows.

$$\eta_i = \begin{cases} \frac{M}{\delta^*} \psi_i \sum_{j=0}^{\lfloor \frac{\delta^*}{\tau_i} \rfloor} \alpha_m^j & \text{for } 1 \leq i \leq k, \\ \frac{M}{\delta^*} \psi_i \alpha_a \sum_{j=0}^{\lfloor \frac{\delta^*}{\tau_i} \rfloor} j & \text{for } k+1 \leq i \leq l. \end{cases} \quad (38)$$

For (36) to be consistent δ^* has to be equal to $\delta^* = \frac{\log[1/\beta_m]}{\log[\alpha_m]} \min_{1 \leq i \leq k} \tau_i$. Therefore, among the MIMD sessions, only the session with the least RTT will have an equilibrium window size different from 0. The equilibrium window of the other MIMD sessions will go to 0. We can, therefore, consider the case where there is only one MIMD session and several AIMD sessions.

For $k = 1$, from (36) and (37), we obtain

$$\begin{aligned} \delta^* &= \tau_1 \frac{\log[1/\beta_m]}{\log \alpha_m}, \\ \psi_i &= \alpha_a \frac{\beta_a}{1 - \beta_a} \frac{\log[1/\beta_a]}{\log \alpha_a} \frac{\tau_1}{\tau_i}, \quad 2 \leq i \leq l, \\ \psi_1 &= \beta_m \Lambda - \frac{\beta_m}{\beta_a} \sum_{i=2}^l l \psi_i. \end{aligned}$$

The inter-loss time depends entirely on the parameters and the RTT of the MIMD session. The effect of different RTTs for the AIMD sessions is to scale ψ_i by a factor of $\frac{\tau_1}{\tau_i}$. Therefore, an AIMD session with lower RTT can obtain a better throughput.

6.1 Several MIMD

If there are l MIMD sessions with different RTTs sharing a link, then the session with the smallest RTT will get all the capacity and the windows for other sessions will go to 0. This result was also mentioned, for two sessions, in [7]. However, if the sessions have the liberty to choose their increase and decrease parameters then each session can obtain some share of the capacity. Let α_{mi} and β_{mi} be the increase parameter and the decrease parameter, respectively, of the i^{th} MIMD session.

Proposition 6.1 *l sessions with different RTTs will have a behaviour similar to l sessions with the same RTT if*

$$\tau_i \frac{\log[1/\beta_{mi}]}{\log[\alpha_{mi}]} = a \quad (a \text{ constant}). \quad (39)$$

The inter-loss time, δ^ , will then be equal to a .*

Proof 6.1 *We note that with this value of δ^* , (36) is consistent. Therefore, an equilibrium solution exists. Let $x(0)$ is the initial window vector. The time to the first synchronized loss, t_1 , can be computed using the condition $\sum_{i=0}^l x_i(0) \alpha_{mi}^{t_1/\tau_i} = \Lambda$. We can now compute $x_i(1) = \beta_{mi} x_i(0) \alpha_{mi}^{t_1/\tau_i}$. The next loss will occur after a time δ^* . This can be verified by noting that $\sum_{i=0}^l \frac{x_i(1)}{\beta_{mi}} = \Lambda$, and $t_2 = \delta^*$ given by (39) satisfies $\sum_{i=0}^l x_i(1) \alpha_{mi}^{t_2/\tau_i} = \Lambda$. From this we obtain $t_2 = \delta^*$. Since t_n is the equilibrium value of the inter-loss time, $x_i(1)$ will also be the equilibrium value of ψ_i . Now, the system will be similar to the same RTT case where the equilibrium window vector is the same as the window vector just after the first synchronized loss.*

Therefore we can obtain a condition on setting the increase and decrease parameters of MIMD algorithm as a function the RTT in order not to be extremely unfair.

7 Conclusions

In the first part of the paper, we studied the fairness in sessions (with either the same RTT or different RTTs) using MIMD congestion control algorithm. For the sessions with the same RTT, it was observed that there was extreme unfairness even in the case of rate independent asynchronous losses. It was shown that fair sharing could be achieved by introducing a stream of rate dependent asynchronous losses. For the case of sessions with different RTTs, it was observed that the arrival rate of these rate dependent asynchronous losses had to be greater than a certain minimum rate in order to achieve fairness. Therefore, in networks with sessions using MIMD algorithms, a stream of rate dependent asynchronous losses, using, for example, some buffer management scheme, would be necessary to ensure fair sharing. In the second part of the paper, we studied capacity sharing between MIMD sessions and AIMD sessions. It was noted that for a given set of parameters, AIMD sessions got a throughput which is independent of the BDP, and that the rest of the capacity was utilized by the MIMD session. In networks with BDP less than a threshold value, it was observed that one AIMD session obtained better throughput than one MIMD session. It was also observed that an AIMD user could open multiple sessions in order to improve its observed throughput whereas for the MIMD user the throughput was invariant to the number of sessions it opened.

References

- [1] T. Kelly, "Scalable TCP: Improving Performance in Highspeed Wide Area Networks," *Computer Communication Review*, vol. 33, no. 2, pp. 83–91, April 2003.
- [2] G. Vinnicombe, "On the Stability of Networks Operating TCP-like Congestion Control," in *Proc. of the IFAC World Congress*, 2002.
- [3] D. Chiu and R. Jain, "Analysis of the Increase/Decrease Algorithms for Congestion Avoidance in Computer Networks," *Journal of Computer Networks and ISDN*, vol. 17, no. 1, pp. 1–14, 1989.
- [4] S. Gorinsky, "Feedback Modeling in Internet Congestion Control," in *Proc. of NEW2AN 2004*, February 2004. Also see a technical report at <http://www.arl.wustl.edu/~gorinsky/TR2002-39.ps>.
- [5] G. Hasegawa, M. Murata, and H. Miyahara, "Fairness and Stability of congestion control mechanisms of TCP," *Telecommunication Systems*, November 2000.

-
- [6] D. Loguinov and H. Radha, “End-to-End Rate-Based Congestion Control: Convergence Properties and Scalability Analysis,” *IEEE/ACM Transactions on Networking*, vol. 11, no. 4, August 2003.
 - [7] L. Xu, K. Harfoush, and I. Rhee, “Binary Increase Congestion Control (BIC) for Fast Long-Distance Network,” in *Proc. of IEEE Infocom*, March 2004.
 - [8] A. Budhiraja, F. Hernández-Campos, V.G. Kulkarni, , and F. D. Smith, “Stochastic Differential Equation for TCP Window Size: Analysis and Experimental Validation,” *Probability in the Engineering and Informational Sciences*, vol. 18, pp. 111–140, 2004.
 - [9] O. Ait-Hellal, E. Altman, D. Elouadghiri, M. Erramdani, and N. Mikou, “Performance of TCP/IP: the case of two Controlled Sources,” in *Proc. of the ICC’97*, 1997, pp. 469–477.
 - [10] G. Fayolle, V.A. Malyshev, and M.V. Menshikov, *Topics in the Constructive Theory of Countable Markov Chains*, Cambridge University Press, 1995.
 - [11] L. Kleinrock, *Queueing Systems Volume I: Theory*, Wiley & sons, 1975.
 - [12] J.G. Kemeny, J.L. Snell, and A.W. Knapp, *Denumerable Markov Chains*, Springer-Verlag, New York, 2nd edition, 1976.
 - [13] S.P. Meyn and R. Tweedie, *Markov Chains and Stochastic Stability*, Springer, London, 1993.
 - [14] S. McCanne and S. Floyd, “ns: Network Simulator,” Available at <http://www.isi.edu/nsnam/ns/>.



Unité de recherche INRIA Sophia Antipolis
2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

Unité de recherche INRIA Futurs : Parc Club Orsay Université - ZAC des Vignes
4, rue Jacques Monod - 91893 ORSAY Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Unité de recherche INRIA Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

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