ELSEVIER

Brief paper

# Modelling TCP congestion control dynamics in drop-tail environments ☆

Robert Shorten [a,*,1], Chris King [b,**,1,2], Fabian Wirth [c], Douglas Leith [a]

[a]*Hamilton Institute, National University of Ireland, Maynooth, Ireland*
[b]*Department of Mathematics, Northeastern University, Boston MA 02115, USA*
[c]*Department of Mathematics, University of Bremen, Germany*

## Abstract

In this paper we study communication networks that employ drop-tail queueing and *additive-increase multiplicative-decrease* (AIMD) congestion control algorithms. We show that the theory of non-negative matrices may be employed to model such networks and to derive basic theorems concerning their behaviour.
© 2007 Published by Elsevier Ltd.

*Keywords:* TCP; Network congestion control

## 1. Introduction

In this paper we describe a modelling approach that captures the essential features of communication networks of additive-increase multiplicative-decrease (AIMD) sources that employ drop-tail and other queuing disciplines. The novelty of our approach lies in the fact that we are able to use the theory of non-negative matrices and hybrid systems to build mathematical models of unsynchronised communication networks by extending the approach first developed in Shorten, Leith, Foy, and Kilduff (2005). In particular, we show that it is possible to relate important network properties to the characteristics of the non-negative matrices that arise in the study of such communication networks under very general conditions.

While an extensive literature exists relating to the modelling of transmission control protocol (TCP) traffic, the models presented in this paper represent a departure from traditional network models. Many recent models are based on the so-called fluid approaches and focus on active queueing disciplines, see for example Srikant (2004) and the references therein. While such models are powerful for design, the applicability of such models for networks with drop-tail buffers, and for networks with low numbers of flows, remain open questions. The well-known square-root formula of Padhye, Firoiu, Towsley, and Kurose (2000) provides an approximate expression for the congestion window achieved by a TCP flow operating in a bath of noise. The statistical independence assumptions in this model however neglect interactions between competing flows, and consequently the dynamics of networks in which TCP operates. Recently several authors have developed new types of models suited to drop-tail networks: most notably by Hespanha (2004) and Baccelli and Hong (2002). We note that while the model derived in Baccelli and Hong (2002) is similar to the model presented in this paper, the work by Baccelli and Hong does not exploit the non-negativity that is central to the work presented here.

Our paper is structured as follows. In Section 2 we develop a positive systems network model that captures the essential features of communication networks employing AIMD congestion control algorithms. This approach gives rise to a model of AIMD networks in which the network dynamics are described by a finite set of non-negative matrices. The main results of this paper are then presented in Section 3. Our main proofs are given in the Appendix.

---

* Corresponding author.
**Corresponding author.
*E-mail address:* robert.shorten@may.ie (R. Shorten).
[1] Joint first authors.
[2] Permanent address: Department of Mathematics, Northeastern University, Boston MA 02115, USA.

## 2. Non-negative matrices and communication networks

*Communication networks*: A communication network consists of a number of sources and sinks connected together via links and routers. We assume that these links can be modelled as a constant propagation delay together with a queue, that the queue is operating according to a drop-tail discipline, and that all of the sources are operating a TCP-like congestion control algorithm.

*AIMD algorithm*: In the original paper proposing the AIMD algorithm, Chiu and Jain (1989) consider a system in which $n$ users compete for a resource having limited availability per unit time, e.g., bandwidth in communication networks. The users' actions consist of a continual gentle probing for the availability of the resource by submitting requests for its use—these requests are satisfied whenever global capacity is not exceeded. Specifically, the probing action consists of additively increasing the send rate according to some rule. The situation is depicted in Fig. 1, with $w_i(t)$ representing the number of units of the resource that user $i = 1, \dots, n$ tries to use at time $t \geqslant 0$. A key assumption in the model formulated by Chiu and Jain is the assertion that the users do not communicate directly with each other. Further, the only information about availability of the resource that the users get is when the collective utilisation of the resource exceeds some capacity constraint. At such time-instances, referred to as *congestion events*, some, or all users are instantly and simultaneously informed through a binary feedback. The users then respond to these notifications of congestion by decentralised down-scaling of their individual utilisation-rates in a multiplicative fashion $w_i(t) \rightarrow \beta_i w_i(t)$ where $\beta_i \in (0, 1)$.

*TCP*: TCP operates a window-based congestion control algorithm that uses the AIMD algorithm to allocate 11 bandwidth between competing network users, the TCP standard defines a variable *cwnd* called the congestion window. Each source uses this variable to track the number of sent unacknowledged packets that can be in transit at any time. When the window size is exhausted, the source must wait for an acknowledgement before sending a new packet. Congestion control is achieved by dynamically adapting the window size according to the AIMD law.
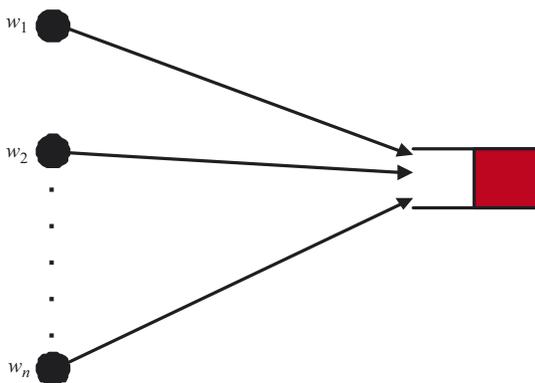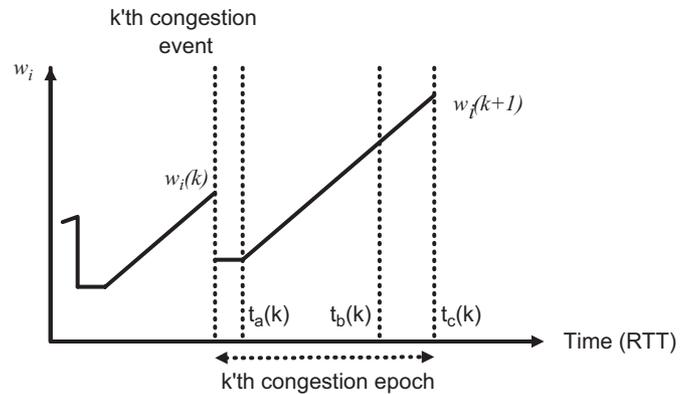


Fig. 1. *n*-player system.



Fig. 2. Evolution of window size.

### 2.1. Synchronised communication networks employing AIMD

For convenience we recall the following discussion from Shorten et al. (2005). The material in the present paper constitutes an extension of Shorten et al. (2005) to unsynchronised networks and it is useful to repeat the following discussion as it will considerably aid the exposition in later sections.

Consider communication networks for which the following assumptions are valid: (i) at congestion every source experiences a packet drop; and (ii) each source has the same round-trip-time (RTT).[3] In this case an exact model of the network dynamics may be found as follows (Shorten et al., 2005).

Let $w_i(k)$ denote the congestion window size of source $i$ immediately before the $k$th network congestion event is detected by the source. Over the $k$th congestion epoch three important events can be discerned: $t_a(k)$, $t_b(k)$ and $t_c(k)$; as depicted in Fig. 2. The time $t_a(k)$ denotes the instant at which the number of unacknowledged packets in flight equals $\beta_i w_i(k)$ where $\beta_i$ is the multiplicative factor of the $i$th flow ($\alpha_i$ is the additive increase factor of this flow); $t_b(k)$ is the time at which the bottleneck queue is full; and $t_c(k)$ is the time at which packet drop is detected by the sources, where time is measured in units of RTT.[4] It follows from the definition of the AIMD algorithm that the window evolution is completely defined over all time instants by knowledge of the $w_i(k)$ and the event times $t_a(k)$, $t_b(k)$ and $t_c(k)$ of each congestion epoch. We therefore only need to investigate the behaviour of these quantities.

We assume that each source is informed of congestion one RTT after the queue at the bottleneck link becomes full; that is $t_c(k) - t_b(k) = 1$. Also when the sources detect congestion, the total window size has reached the capacity $P$ of the pipe and each source has increased its window size one more time

---

[3] One RTT is the time between sending a packet and receiving the corresponding acknowledgement when there are no packet drops.

[4] Note that measuring time in units of RTT results in a linear rate of increase for each of the congestion window variables between congestion events.

before packet loss is detected. That is

$$w_i(k) \geqslant 0 \quad \text{and} \quad \sum_{i=1}^{n} w_i(k) = P + \sum_{i=1}^{n} \alpha_i \ \forall k > 0, \tag{1}$$

where $P$ is the maximum number of packets which can be in transit in the network at any time; $P$ is usually equal to $q_{max} + BT_d$ where $q_{max}$ is the maximum queue length of the congested link, $B$ is the service rate of the congested link in packets per second and $T_d$ is the round-trip time when the queue is empty. At the $(k+1)$th congestion event

$$w_i(k+1) = \beta_i w_i(k) + \alpha_i[t_c(k) - t_a(k)], \tag{2}$$

and summing over all sources yields

$$
\begin{aligned}
&t_c(k) - t_a(k) \\
&= \frac{1}{\sum_{i=1}^{n} \alpha_i} \left[ P - \sum_{i=1}^{n} \beta_i w_i(k) \right] + 1.
\end{aligned}
\tag{3}
$$

Hence, it follows that

$$w_i(k+1) = \beta_i w_i(k) + \frac{\alpha_i}{\sum_{j=1}^{n} \alpha_j} \left[ \sum_{j=1}^{n} (1 - \beta_j) w_j(k) \right] \tag{4}$$

and so the dynamics of an entire network of such sources is given by

$$W(k+1) = AW(k), \tag{5}$$

where $W^{\mathrm{T}}(k) = [w_1(k), \ldots, w_n(k)]$, and

$$
A = 
\begin{bmatrix}
\beta_1 & 0 & \cdots & 0 \\
0 & \beta_2 & 0 & 0 \\
\vdots & 0 & \ddots & 0 \\
0 & 0 & \cdots & \beta_n
\end{bmatrix}
+ \frac{1}{\sum_{j=1}^{n} \alpha_j}
$$

$$
\times
\begin{bmatrix}
\alpha_1 \\
\alpha_2 \\
\cdots \\
\alpha_n
\end{bmatrix}
[1 - \beta_1 \ \ 1 - \beta_2 \ \ \cdots \ \ 1 - \beta_n]. \tag{6}
$$

The matrix $A$ is a positive matrix (all the entries are positive real numbers) and it follows that the synchronised network (5) is a positive linear system (Berman & Plemmons, 1979). Many results are known for positive matrices and we exploit some of these to characterise the properties of synchronised communication networks. In particular, from the viewpoint of designing communication networks the following properties are important: (i) network fairness; (ii) network convergence and responsiveness; and (iii) network throughput. It is shown in Shorten et al. (2005) that these properties can be deduced from the network matrix A. In particular:

**Theorem 2.1** (*Shorten et al., 2005*). *Let A be defined as in Eq.* (6). *Then A is a column stochastic matrix with Perron eigenvector $x_p^{\mathrm{T}} = [\alpha_1/(1 - \beta), \ldots, \alpha_n/(1 - \beta_n)]$ and whose eigenvalues are real and positive. Further, the network converges to a unique stationary point $W_{ss} = \Theta x_p$, where $\Theta$ is*

*a positive constant such that the constraint* (1) *is satisfied;* $\lim_{k \to \infty} W(k) = W_{ss}$; *and the rate of convergence of the network to $W_{ss}$ is bounded by the second largest eigenvalue of A.*

### 2.2. Unsynchronised networks

To distinguish variables in this section, we denote the nominal parameters of the sources used in the previous section by $\alpha_i^s, \beta_i^s, i = 1, \ldots, n$. Now consider the general case of a number of sources competing for shared bandwidth in a generic dumbbell topology. As before, in our discussion $k$ is still used to enumerate congestion epochs. Note, however, that a congestion epoch is the time elapsing between one instant when packets are lost by some source and the next instant when packets are lost *by the same or some other source*. That is, congestion epochs are now globally defined with respect to the bottleneck router. As before a number of important events may be discerned, where we now measure time in seconds, rather than units of RTT. Denote by $t_{ai}(k)$ the time at which the number of packets in flight belonging to source $i$ is equal to $\beta_i^s w_i(k)$; $t_q(k)$ is the time at which the bottleneck queue begins to fill; $t_b(k)$ is the time at which the bottleneck queue is full; and $t_{ci}(k)$ is the time at which the $i$th source is informed of congestion. In this case the evolution of the $i$th congestion window variable does not evolve linearly with time after $t_q$ seconds due to the effect of the bottleneck queue filling and the resulting variation in RTT. Note also that we do not assume that every source experiences a drop when congestion occurs.

Given these general features it is clear that the modelling task is more involved than in the synchronised case. Nonetheless, it is possible to relate $w_i(k)$ and $w_i(k+1)$ using the same approach as before.

Specifically, we now allow the $i$th source to experience congestion at the end of the epoch whereas the $j$th source does not. This corresponds to the $i$th source reducing its congestion window by the factor $\beta_i^s$ after the $k+1$th congestion event, and the $j$th source not adjusting its window size at the congestion event. We therefore allow the back-off factor of the $i$th source to take one of two values at the $k$th congestion event

$$\beta_i(k) \in \{\beta_i^s, 1\}, \tag{7}$$

corresponding to whether the source experiences a packet loss or not.

Due to the variation in round trip time, the congestion window of a flow does not evolve linearly with time over a congestion epoch. Nevertheless, we may relate $w_i(k)$ and $w_i(k+1)$ linearly by defining an average rate $\alpha_i(k)$ depending on the $k$th congestion epoch

$$\alpha_i(k) := \frac{w_i(k+1) - \beta_i(k)w(k)}{T(k)}, \tag{8}$$

where $T(k)$ is the duration of the $k$th epoch measured in seconds. Equivalently we have

$$w_i(k+1) = \beta_i(k)w_i(k) + \alpha_i(k)T(k). \tag{9}$$

Recall, that $T_{d_i}$ denotes the round trip time of source $i$, when the queue is empty and $R_{d_i} := T_{d_i} + q_{max}/B$ is the RTT,

when the queue is full. The number of round trips of source $i$ during the time interval $T(k)$ is in the interval $[T(k)/R_{d_i}, T(k)/T_{d_i}]$ and therefore

$$\alpha_i(k) \in \left[\frac{\alpha_i}{R_{d_i}}, \frac{\alpha_i}{T_{d_i}}\right].$$

In the case when $q_{\max} \ll BT_{d_i}$, $i = 1, \ldots, n$, we have $T_{d_i} \approx R_{d_i}$ and the average $\alpha_i$ are (almost) independent of $k$ and given by $\alpha_i(k) \approx \alpha_i^s/T_{d_i}$ for all $k \in \mathbb{N}$, $i = 1, \ldots, n$.

The approximation

$$\alpha_i \approx \frac{\alpha_i^s}{T_{d_i}}, \quad i = 1, \ldots, n \tag{10}$$

is of considerable practical importance and corresponds to the case of a network whose bottleneck buffer is small compared with the delay-bandwidth product. In view of (7) and (9) a convenient representation of the network dynamics is obtained as follows. At congestion the bottleneck link is operating at its capacity $B$, i.e.,

$$\sum_{i=1}^{n} \frac{w_i(k) - \alpha_i}{RTT_{i,\max}} = B, \tag{11}$$

where $RTT_{i,\max}$ is the RTT experienced by the $i$th flow when the bottleneck queue is full. Note, that $RTT_{i,\max}$ is independent of $k$. Setting $\gamma_i := (RTT_{i,\max})^{-1}$ we have that

$$\sum_{i=1}^{n} \gamma_i w_i(k) = B + \sum_{i=1}^{n} \gamma_i \alpha_i. \tag{12}$$

By interpreting (12) at $k + 1$ and inserting (9) for $w_i(k+1)$ it follows that

$$\sum_{i=1}^{n} \gamma_i \beta_i(k) w_i(k) + \gamma_i \alpha_i T(k) = B + \sum_{i=1}^{n} \gamma_i \alpha_i. \tag{13}$$

Using (12) again it follows that

$$T(k) = \frac{1}{\sum_{i=1}^{n} \gamma_i \alpha_i} \left(\sum_{i=1}^{n} \gamma_i (1 - \beta_i(k)) w_i(k)\right). \tag{14}$$

Inserting this expression into (9) and considering all sources, the dynamics of the entire network of sources at the $k$th congestion event are described by

$$W(k+1) = A(k)W(k), \quad A(k) \in \{A_1, \ldots, A_m\}, \tag{15}$$

where

$$A(k) = \begin{bmatrix} \beta_1(k) & 0 & \cdots & 0 \\ 0 & \beta_2(k) & 0 & 0 \\ \vdots & 0 & \ddots & 0 \\ 0 & 0 & \cdots & \beta_n(k) \end{bmatrix} + \frac{1}{\sum_{j=1}^{n} \gamma_j \alpha_j}$$
$$\times \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \cdots \\ \alpha_n \end{bmatrix} [\gamma_1(1 - \beta_1(k)), \ldots, \gamma_n(1 - \beta_n(k))]$$

and where $\beta_i(k)$ is either 1 or $\beta_i^s$. The nonnegative matrices $A_2, \ldots, A_m$ are constructed by taking the matrix $A_1$,

$$A_1 = \begin{bmatrix} \beta_1^s & 0 & \cdots & 0 \\ 0 & \beta_2^s & 0 & 0 \\ \vdots & 0 & \ddots & 0 \\ 0 & 0 & \cdots & \beta_n^s \end{bmatrix} + \frac{1}{\sum_{j=1}^{n} \gamma_j \alpha_j}$$
$$\times \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \cdots \\ \alpha_n \end{bmatrix} [\gamma_1(1 - \beta_1^s), \ldots, \gamma_n(1 - \beta_n^s)]$$

and setting some, but not all, of the $\beta_i$ to 1. This gives rise to $m = 2^n - 1$ matrices associated with the system (15) that correspond to the different combinations of source drops that are possible. We denote the set of these matrices by $\mathscr{A}$.

*Comment* 1: Note that another, sometimes very useful, representation of the network dynamics can be obtained by considering the evolution of scaled window sizes at congestion; namely, by considering the evolution of the vectors of network states $W_\gamma^{\mathrm{T}}(k) = [\gamma_1 w_1(k), \gamma_2 w_2(k), \ldots, \gamma_n w_n(k)]$. Here one obtains the following description of the network dynamics:

$$W_\gamma(k+1) = \overline{A}(k) W_\gamma(k) \tag{16}$$

with $\overline{A}(k) \in \overline{\mathscr{A}} = \{\overline{A}_1, \ldots, \overline{A}_m\}$, $m = 2^n - 1$ and where the $\overline{A}_i$ are obtained by the similarity transformation associated with the change of variables. All of the matrices in the set $\overline{\mathscr{A}}$ are now column stochastic as in the synchronised case.

*Comment* 2: Our discussion extends the (exact) synchronised model to the unsynchronised case. To do this, we have effectively assumed linear probing action between congestion events. This assumption allows us to develop nonnegative matrix models of general AIMD networks with general linear capacity constraints $C = \Theta^{\mathrm{T}} W(k)$, where $\Theta$ defines the normal to a linear hyperplane:

$$W(k+1) = A(k)W(k), \quad A(k) \in \{A_1, \ldots, A_m\}, \tag{17}$$

with $A(k) = B(k) + \alpha\Theta^{\mathrm{T}}/\Theta^{\mathrm{T}}\alpha(I - B(k))$, where $\alpha^{\mathrm{T}} = [\alpha_1, \ldots, \alpha_n]$ and where the matrix $B(k)$ is given by

$$B(k) = \begin{bmatrix} \beta_1(k) & 0 & \cdots & 0 \\ 0 & \beta_2(k) & 0 & 0 \\ \vdots & 0 & \ddots & 0 \\ 0 & 0 & \cdots & \beta_n(k) \end{bmatrix}.$$

## 3. Mathematical results

It follows from (15) that $W(k) = \Pi(k)W(0)$, where $\Pi(k) = A(k)A(k-1) \ldots A(0)$. Consequently, the behaviour of $W(k)$, as well as the network fairness and convergence properties, are governed by the properties of the infinite matrix product $\Pi(k)$. The objective of this section is to analyse the average behaviour of $\Pi(k)$ with a view to making concrete statements about important network properties.

### 3.1. Constant drop-probabilities

Here, we make two simplifying assumptions.

**Assumption 3.1.** The probability that $A(k) = A_i$ in (15) is independent of $k$ and equals $\rho_i$.

In other words Assumption 3.1 says that the probability that the network dynamics are described by $W(k + 1) = A(k)W(k)$, $A(k) = A_i$ over the $k$th congestion epoch is $\rho_i$ and that the random variables $A(k)$, $k \in \mathbb{N}$ are independent and identically distributed (i.i.d.).

Given the probabilities $\rho_i$ for $i \in \{1, \dots, 2^n - 1\}$, one may then define the probability $\Lambda_j$ that source $j$ experiences a back-off at the $k$th congestion event as follows:

$$\Lambda_j = \sum \rho_i,$$

where the summation is taken over those $i$ which correspond to a matrix in which the $j$th source sees a drop. Or to put it another way, the summation is over those indices $i$ for which the matrix $A_i$ is defined with a value of $\beta_j \neq 1$.

**Assumption 3.2.** We assume that $\Lambda_j > 0$ for all $j \in \{1, \dots, n\}$.

Simply stated, Assumption 3.2 states that almost surely all flows must see a drop at some time (provided that they persist for a long enough time). A consequence of the above assumptions is that the probability that source $j$ experiences a drop at the $k$th congestion event is not independent of the other sources. For example, if the first $n - 1$ sources do not see a drop then this implies that source $n$ must see a drop (in accordance with the usual notion of a congestion event, we require at least one flow to see a drop at each congestion event). Hence, the events cannot be independent.

Under the foregoing assumptions we have the following key result.

**Theorem 3.1.** *Consider the stochastic system defined in the above preamble. Let $\Pi(k)$ be the random matrix product arising from the evolution of the first $k$ steps of this system:*

$$\Pi(k) = A(k)A(k - 1) \dots, A(0).$$

*Then, the expectation of $\Pi(k)$ is given by*

$$E(\Pi(k)) = \left( \sum_{i=1}^{m} \rho_i A_i \right)^k; \tag{18}$$

*and the asymptotic behaviour of $E(\Pi(k))$ satisfies*

$$\lim_{k \to \infty} E(\Pi(k)) = x_p y_p^{\mathrm{T}}, \tag{19}$$

*where the vector $x_p$ is given by $x_p^{\mathrm{T}} = \Theta(\alpha_1/\Lambda_1(1 - \beta_1), \alpha_2/\Lambda_2(1 - \beta_2), \dots, \alpha_n/\Lambda_n(1 - \beta_n))$, $y_p^{\mathrm{T}} = (\gamma_1, \dots, \gamma_n)$. Here $\Theta \in \mathbb{R}$ is chosen such that Eq. (12) is satisfied if $w_i$ is replaced by $x_{pi} = \Theta\alpha_i/(\Lambda_i(1 - \beta_i))$.*

**Corollary 3.1.** *For given $W(0)$ define random variable $\overline{W}(k)$ with*

$$\overline{W}(k) := \frac{1}{k + 1} \sum_{i=0}^{k} W(i).$$

*Then expectation of $\overline{W}(k)$ is given by*

$$E(\overline{W}(k)) = \frac{1}{k + 1}(I + E(A(1)) + E(A(1))2$$
$$+ \dots + E(A(1))^k)W(0).$$

*And since $E(A(1))^k \to x_p y^{\mathrm{T}}$ as $k \to \infty$,*

$$\lim_{k \to \infty} E(\overline{W}(k)) = x_p y^{\mathrm{T}} W(0).$$

The following facts follow immediately from Theorem 3.1.

(*i*) *Convergence*: The congestion window vector $W(k)$ converges, on average, to the unique value $W_{ss} = \Theta x_p$ where $\Theta$ is a positive constant such that the constraint (12) is satisfied. When the $\Lambda_i$, $i = 1, \dots, n$ are equal, $x_p$ is identical to the Perron eigenvector obtained in the case of synchronised networks; that is, the ensemble average in the unsynchronised case is identical to the fixed point in the deterministic situation where packet drops are synchronised.

(*ii*) *Fairness*: Window fairness is achieved, on average, when the vector $x_p$ is a scalar multiple of the vector $[1, \dots, 1]$; that is, when the ratio $\alpha_i/\Lambda_i(1 - \beta_i)$ does not depend on $i$. Observe that unlike in the synchronised case, fairness now depends on the relative drop probability of each flow. When the flows have equal drop probability $\Lambda_i$ then the foregoing fairness condition is identical to that in the synchronised case.

(*iii*) *Network responsiveness*: The magnitude of the second largest eigenvalue $\lambda_2$ of the matrix $\sum_{i=1}^{m} \rho_i A_i$ bounds the convergence properties of the network.

### 3.2. Place dependent drop-probabilities

We now relax the assumption that flows are dropped with constant probabilities at congestion events. Instead, we allow the drop probabilities to depend on the current state of the system, that is the set of current window sizes $W(k)$ at the $k$th congestion event. As an example, it is reasonable to expect at a congestion event that flows with larger window sizes are more likely to be dropped, and this can be realised in our model by making the drop probabilities increasing functions of the current window size.

Mathematically, our model is a discrete time Markov chain whose state space is the simplex $\mathscr{S} = \{w = (w_1, \dots, w_n) : w_i \geqslant 0, \sum w_i = C\}$ where $C$ is the link capacity. The state of the system $W(k)$ evolves according to the rule $W(k + 1) = A(k)W(k)$, where $A(k) \in \mathscr{A}$ is chosen randomly using a place-dependent probability distribution on $\mathscr{A}$. Specifically, for each $w \in \mathscr{S}$ there is a probability distribution $\{p_1(w), p_2(w), \dots, p_m(w)\}$ on $\mathscr{A}$, and $A(k)$ is chosen randomly with probability $P(A(k) = A_i) = p_i(W(k))$. With the

following mild assumptions on the drop probability functions we will show that this model is ergodic.

**Assumption 3.3.** The distribution $p_i(w)$ is uniformly Lipschitz on $\mathscr{S}$ with respect to the $l_1$-norm. That is, we assume there is a constant $K$ such that for all $w, v \in \mathscr{S}$ and all $i = 1, \ldots, m$,

$$|p_i(w) - p_i(v)| \leqslant K\|w - v\|_1 = K\sum_{j=1}^{n}|w_j - v_j|. \tag{20}$$

To set up the notation for the statement of the second assumption, let $s = (i_1, i_2, \ldots, i_M)$ be a $M$-string of indices, labelling a sequence of congestion matrices $A_{i_1}, A_{i_2}, \ldots$ at $M$ successive congestion events, and define the corresponding matrix product $\prod_M(s) = A_{i_M} A_{i_{M-1}} \ldots A_{i_1}$. Given an initial vector $w \in \mathscr{S}$, the Markov chain defines a probability distribution $p^*(\cdot; w)$ on these $M$-strings by

$$p^*(s; w) = p_{i_1}(w) p_{i_2}(A_{i_1}w) \ldots p_{i_M}(A_{i_M} \cdots A_{i_1}w). \tag{21}$$

As before we let $A_1$ denote the matrix corresponding to a congestion event where all flows experience a drop.

**Assumption 3.4.** (A) There is $q > 0$ and a subset $\mathscr{H} \subset \mathscr{S}$, such that $\mathscr{H}$ is mapped into itself by $A_1$, and $p_1(w) \geqslant q$ for every $w \in \mathscr{H}$.

(B) There is an integer $N \geqslant 1$, and $q' > 0$, such that for any $v, w \in \mathscr{S}$, there is a $(N-1)$-string $s$ for which $p^*(s; w) \geqslant q'$, and $\prod_{N-1}(s)w, \prod_{N-1}(s)v \in \mathscr{H}$.

The following theorem extends results in Wirth, Stanojevic, Shorten, and Leith (2006) and presents an alternative treatment of the results in Leizarowitz, Stanojevic, and Shorten (2006).

**Theorem 3.2.** *Assume that the place-dependent drop probabilities $\{p_i(w)\}$ satisfy Assumptions 3.3 and 3.4. Then (i) there is an attractive, unique stationary probability measure for the Markov process $\{W(k)\}$; (ii) for any continuous function $f(w)$ on $\mathscr{S}$, the conditional expectations, $\mathrm{E}[f(W(k))|W(0) = w]$ converge uniformly to a constant as $k \to \infty$; (iii) for any continuous function $f(w)$ on $\mathscr{S}$, the time average $(1/K)\sum_{k=1}^{K} f(W(k))$ converges almost surely to the ensemble average of $f(W)$ with respect to the stationary measure.*

We will prove Theorem 3.2 in the Appendix using results of Isaac (1962), Barnsley, Demko, Elton, and Geronimo (1988), and Stenflo (2002), who established general conditions for ergodicity of Markov chains with place-dependent probabilities. Theorem 3.2 implies that the process $\{W(k)\}$ converges almost surely as $k \to \infty$, and that the limiting distribution is independent of the initial conditions. This ergodic property allows us to relate time averages to ensemble averages, and hence to use pathwise calculations to compute average quantities. Using this method we will show that a version of the result (19) derived in the case of constant drop-probabilities continues to hold for the place-dependent model. Our result will involve the average window size for the $i$th flow computed only at the congestion events where it experiences a drop. To set up the notation, let

$D(k) \subset \{1, \ldots, m\}$ denote the set of flows which experience a drop at the $k$th congestion event, and define

$$\theta_i(k) = 1 \quad \text{if } i \in D(k), \tag{22}$$
$$= 0 \quad \text{if } i \notin D(k). \tag{23}$$

**Theorem 3.3.** *Under the conditions for Theorem 3.2, the following limits exist and are independent of initial conditions:*

$$\langle w_i \rangle = \lim_{k \to \infty} \mathrm{E}[w_i(k)|i \in D(k)], \quad \Lambda_i = \lim_{k \to \infty} \mathrm{E}[\theta_i(k)]. \tag{24}$$

*Furthermore these quantities are related by*

$$\langle w_i \rangle = \frac{\alpha_i}{\Lambda_i(1 - \beta_i)} \mathrm{E}[T], \tag{25}$$

*where $\mathrm{E}[T]$ is the average time between congestion events.*

**Proof.** For any $k \geqslant 1$

$$\mathrm{E}[w_i(k)|i \in D(k)] = \frac{\mathrm{E}[w_i(k)\theta_i(k)]}{\mathrm{E}[\theta_i(k)]}. \tag{26}$$

Furthermore

$$\mathrm{E}[w_i(k)\theta_i(k)] = \mathrm{E}[w_i(k) \sum_{j \in \Lambda_i} p_j(W(k))], \tag{27}$$

where $\Lambda_i \subset \{1, \ldots, m\}$ is the list of all subsets for which the $i$th flow experiences a drop. Applying Theorem 3.2 we conclude that (27) converges uniformly to a value independent of initial conditions as $k \to \infty$. The same argument applies to the denominator in (26), hence the left side of (26) converges to a limit which we define to be $\langle w_i \rangle$. Similar reasoning applies to define $\Lambda_i$.

Considering a sample path of the process, we see that the $i$th window sizes $w_i, w_i'$ at two successive events where flow $i$ is dropped are related by

$$w_i' = \beta_i w_i + \alpha_i T_i \tag{28}$$

where $T_i$ is the time between these events. Define the long-run time averages

$$\langle w_i \rangle_K = \frac{\sum_{k=1}^{K} w_i(k)\theta_i(k)}{\sum_{k=1}^{K} \theta_i(k)}, \quad \langle T_i \rangle_K = \frac{\sum_{k=1}^{K} T(k)}{\sum_{k=1}^{K} \theta_i(k)}, \tag{29}$$

where $T(k)$ is the time between the $k$th and $(k+1)$th congestion events. Then (28) implies

$$\langle w_i \rangle_K = \frac{\alpha_i}{1 - \beta_i} \langle T_i \rangle_K + O\left(\frac{1}{K}\right), \tag{30}$$

where the error term $O(1/K)$ takes care of the mismatches in the sums at $k = 1$ and $K$ (recall that $w_i$ and $T(k)$ are uniformly bounded, so this term is bounded by a constant times $1/K$). Therefore (25) follows from part (iii) of Theorem 3.2, which states that time averages converge to ensemble averages, and hence $\langle w_i \rangle_K$ converges to $\langle w_i \rangle$ and $\langle T_i \rangle_K$ converges to $\mathrm{E}[T]/\Lambda_i$. $\square$

Our next result involves the throughput for the $i$th flow, which is defined by the pathwise expression

$$\delta_i = \lim_{T \to \infty} \frac{1}{T} \int_0^T w_i(t) \, \mathrm{d}t. \tag{31}$$

**Theorem 3.4.** *Under the conditions for Theorem* 3.2, *with probability one the expression* (31) *exists and is independent of the sample path, and is given by*

$$\delta_i = \frac{\Lambda_i(1 - \beta_i^2)}{2\alpha_i} \mathrm{E}[T] \langle w_i^2 \rangle, \tag{32}$$

*where* $\langle w_i^2 \rangle = \lim_{k \to \infty} \mathrm{E}[w_i(k)^2 | i \in D(k)]$. *It satisfies the bounds*

$$\langle w_i \rangle \frac{1 + \beta_i}{2} \leqslant \delta_i \leqslant \langle w_i \rangle \frac{1 + \beta_i}{2} \left( 1 + \frac{\mathrm{VAR}[T]}{\mathrm{E}[T]^2} \right), \tag{33}$$

*where* $\mathrm{VAR}[T]$ *is the variance of the time between congestion events.*

**Proof.** We will write $\{\tau_1, \tau_2, \ldots\}$ to denote the times of the congestion events where flow $i$ experiences a reduction, and $\{w_i(1), w_i(2), \ldots\}$ its window sizes at these events. The evolution equation for $w_i$ between congestion events is

$$w_i(k + 1) = \beta_i w_i(k) + \alpha_i(\tau_{k+1} - \tau_k). \tag{34}$$

Elementary calculations along the sample path show that $\int w_i(t) \, \mathrm{d}t$ can be expressed as a sum of the squares of window sizes at congestion events where the $i$th flow is dropped. Ergodicity then relates this time average to the ensemble average, and this gives (32).

We now use upper and lower bounds on $\langle w_i^2 \rangle$ to derive (33). For the lower bound we just use

$$\langle w_i^2 \rangle \geqslant (\langle w_i \rangle)^2. \tag{35}$$

For the upper bound we square (34) and take the expected value to get

$$\mathrm{E}[w_i(k + 1)^2] = \beta_i^2 \mathrm{E}[w_i(k)^2] + \alpha_i^2 \mathrm{E}[(\tau_{k+1} - \tau_k)^2]$$
$$+ 2\alpha_i \beta_i \mathrm{E}[w_i(k)(\tau_{k+1} - \tau_k)]. \tag{36}$$

We use the bound $2w_i(k)(\tau_{k+1} - \tau_k) \leqslant x w_i(k)^2 + x^{-1}(\tau_{k+1} - \tau_k)^2$, which holds for every $x \geqslant 0$; inserting into (36) gives

$$\mathrm{E}[w_i(k + 1)^2] \leqslant (\beta_i^2 + \alpha_i \beta_i x) \mathrm{E}[w_i(k)^2]$$
$$+ (\alpha_i^2 + \alpha_i \beta_i x^{-1}) \mathrm{E}[(\tau_{k+1} - \tau_k)^2]. \tag{37}$$

Taking $k \to \infty$ and using the ergodic property gives

$$\langle w_i^2 \rangle \leqslant \frac{\alpha_i^2 + \alpha_i \beta_i x^{-1}}{1 - \beta_i^2 - \alpha_i \beta_i x} \mathrm{E}[T^2] \tag{38}$$

where $T$ is the time between congestion events. Using the optimal value $x = \alpha_i^{-1}(1 - \beta_i)$ we find

$$\langle w_i^2 \rangle \leqslant \frac{\alpha_i^2}{(1 - \beta_i)^2} \mathrm{E}[T^2] \tag{39}$$

Finally, we use $\mathrm{E}[T^2] = \mathrm{VAR}[T] + (\mathrm{E}[T])^2$ to get the result. $\square$

### 3.3. General capacity constraints

A network of routers is described by a collection of flows with general capacity constraints of the form $\sum_{i \in R_a} w_i \leqslant C_a$, where $C_a$ is the capacity of the router $R_a$. We consider the case where each router drops flows at its congestion events according to some random protocol, which may or may not depend on the present or past state of the flow window sizes at the router. We can analyse the behaviour of this network of flows by *assuming ergodicity of the process*, and then using pathwise calculations to deduce average properties of the process. Since the pathwise computations are identical to the ones used in the previous case of the single-router Markov model, we find that the same relations (25) hold between ensemble averages in the general case.

## 4. Conclusions

In this paper we have presented a random matrix model that describes the behaviour of a network of $n$ AIMD flows that compete for shared bandwidth via a bottleneck router employing drop-tail queuing. We have used this model to relate several important network properties to properties of sets of nonnegative matrices that arise in the study of such networks. We have also derived under simplifying assumptions a number of analytic results that characterise the ensemble-average throughput of such networks. Finally, we note that we have also validated our results against packet-level simulations for networks many flows and the results are presented in Shorten, Wirth, and Leith (2006).

## Appendix A. Proof of mathematical results

### A.1. Constant probabilities

It was noted before that the matrices in the set $\mathscr{A}$ are not column stochastic. However, the matrices in this set are simultaneously similar to a set of column stochastic matrices under the transformation $\Gamma = \mathrm{diag}[\gamma_1, \ldots, \gamma_n]$. The corresponding transformed dynamics are given by (16) and define the evolution of the vector $\overline{W}_\gamma(k)$. The corresponding results for the system

(15) are directly deduced from these results by similarity. In the following it will be convenient to introduce a notation that identifies each matrix $\overline{A} \in \mathscr{A}$ with the sources that do not see a drop in that congestion event. Let $\mathscr{I} \subset \{1, 2, \ldots, n\}$ be the index set of sources not experiencing congestion at a congestion event. (Clearly, $\mathscr{I} = \{1, 2, \ldots, n\}$ can be ignored, as this means that there is no congestion.) The matrix corresponding to an index set $\mathscr{I}$ is given by

$$\overline{A}_{\mathscr{I}} = \mathrm{diag}(\beta_1(\mathscr{I}), \ldots, \beta_n(\mathscr{I}))$$
$$+ \overline{c}_\alpha \overline{\alpha}[1 - \beta_1(\mathscr{I}), \ldots, 1 - \beta_n(\mathscr{I})],$$

where $\beta_i(\mathscr{I}) = 1$, if $i \in \mathscr{I}$ and $\beta_i(\mathscr{I}) = \beta_i^s$ otherwise and $\overline{c}_\alpha := (\sum_{j=1}^n \gamma_j \alpha_j)^{-1}$ and $\overline{\alpha}^{\mathrm{T}} = [\gamma_1 \alpha_1, \ldots, \gamma_n \alpha_n]$. We now recover our set of possible matrices by

$$\overline{\mathscr{A}} := \{\overline{A}_{\mathscr{I}} | \mathscr{I} \subsetneq \{1, 2, \ldots, n\}\}, \tag{40}$$

which results in a set of $2^n - 1$ matrices, as it should. Note that all $\overline{A} \in \overline{\mathscr{A}}$ are column stochastic, so that they have an eigenvalue equal to 1 equal to the spectral radius. In the following we will use the notation $\overline{A}_{\mathscr{I}} = \Delta_{\mathscr{I}} + \overline{c}_\alpha \overline{\alpha} \beta(\mathscr{I})^{\mathrm{T}}$, where $\Delta_{\mathscr{I}}$ denotes the diagonal matrix and $\beta(\mathscr{I})$ is the vector with entries $1 - \beta_i(\mathscr{I})$. We denote by $\overline{\Pi}(k)$ products of length $k$ of matrices $\overline{A} \in \overline{\mathscr{A}}$.

**Lemma 4.1.** *Consider the random system* (15) *subject to Assumptions* 3.1 *and* 3.2. *The expectation of* $\overline{\Pi}(k)$ *is*

$$E(\overline{\Pi}(k)) = \left( \sum_{\mathscr{I}} \rho_{\mathscr{I}} \overline{A}_{\mathscr{I}} \right)^k. \tag{41}$$

**Proof.** Expanding the power relation on the left-hand side, by independence we have that the expectation of the product is the product of the expectations. This implies the equality. $\quad\square$

**Lemma 4.2.** *Assume that* $\Lambda_i > 0$ *for* $i = 1, \ldots, n$ *then the expectation*

$$E(\overline{A}) = \sum_{\mathscr{I}} \rho_{\mathscr{I}} \overline{A}_{\mathscr{I}}$$

*is positive, column stochastic, and a Perron eigenvector for it is given by*

$$\overline{x}_p^{\mathrm{T}} = \left( \frac{\overline{\alpha}_1}{\Lambda_1(1 - \beta_1)}, \frac{\overline{\alpha}_2}{\Lambda_2(1 - \beta_2)}, \ldots, \frac{\overline{\alpha}_n}{\Lambda_n(1 - \beta_n)} \right). \tag{42}$$

**Proof.** By definition of the expectation we have

$$E(\overline{A}) = \sum_{\mathscr{I}} \rho_{\mathscr{I}} \overline{A}_{\mathscr{I}} = \sum_{\mathscr{I}} \rho_{\mathscr{I}} \Delta_{\mathscr{I}} + \overline{c}_\alpha \sum_{i=1}^m \rho_{\mathscr{I}} \overline{\alpha} \beta(\mathscr{I})^{\mathrm{T}}$$
$$= E(\Delta) + \overline{c}_\alpha \overline{\alpha} E(\beta)^{\mathrm{T}}. \tag{43}$$

The $i$th diagonal entry of the diagonal matrix $E(\Delta)$ is

$$E(\Delta_{i,i}) = \Lambda_i \beta_i + (1 - \Lambda_i) \tag{44}$$

and the $i$th entry of $E(\beta)$ is

$$E(\beta_i) = \Lambda_i(1 - \beta_i). \tag{45}$$

Hence, the matrix $E(\overline{A})$ is of the form of $\overline{A}_1$, with the same vector $\overline{\alpha}$ and where $\beta_i$ replaced by $\tilde{\beta}_i := 1 - \lambda_i(1 - \beta_i) \in (0, 1)$. It follows by Theorem 2.1 that a Perron eigenvector of $E(\overline{A})$ is given by $\overline{x}_p^{\mathrm{T}} = (\overline{\alpha}_1/\Lambda_1(1 - \beta_1), \overline{\alpha}_2/\Lambda_2(1 - \beta_2), \ldots, \overline{\alpha}_n/\Lambda_n(1 - \beta_n))$. $\quad\square$

Theorem 3.1 follows immediately from the above two Lemmas and Theorem 2.1.

### A.2. Place dependent probabilities

**Proof of Theorem 3.2.** Barnsley et al. (1988) have derived a general condition for ergodicity of Markov chains with place-dependent probabilities, and their results are summarised in Stenflo's (2002) paper. The version closest to our model (compact state space, Lipschitz continuous transition functions) was considered earlier by Isaac (1962). Elton (1990) extended these results by proving almost sure convergence of time averages to ensemble averages under the same conditions. The conditions which guarantee ergodicity in our case are contained in the following lemma. These conditions state that the $N$-step transition probabilities satisfy an average contractivity property.

**Lemma 4.3.** *Assume the drop probabilities satisfy the Assumptions* 3.3 *and* 3.4. *Then there is* $r < 1$ *such that for all* $v, w \in \mathscr{S}$,

$$\sum_s p^*(s; w) \|\Pi_N(s)v - \Pi_N(s)w\|_1 \leqslant r\|v - w\|_1. \tag{46}$$

*Furthermore, given any* $\varepsilon > 0$, *there is an integer* $M$, *and* $\delta > 0$, *such that for any* $v, w \in \mathscr{S}$

$$\sum_{s,t : \|\Pi_M(s)v - \Pi_M(t)w\|_1 < \varepsilon} p^*(s; v) p^*(t; w) > \delta. \tag{47}$$

**Proof.** As noted in Comment 1, we can use a similarity transformation to ensure that each matrix $A_i$ is column stochastic. Hence $A_i$ is a contraction with respect to the $l_1$ norm, that is for any vector $x$ and index $i$ we have

$$\|A_i x\|_1 \leqslant \|x\|_1. \tag{48}$$

Also, recall that $A_1$ is the matrix corresponding to the case where all flows are dropped. Then $A_1$ is an entrywise positive column stochastic matrix, and hence is a strict contraction with respect to the $l_1$ norm, so there is $r' < 1$ such that

$$\|A_1 x\|_1 \leqslant r'\|x\|_1. \tag{49}$$

From (48) we get $\|\Pi_M(s)v - \Pi_M(s)w\|_1 \leqslant \|v - w\|_1$ for all $M$-strings $s$, all $M$ and all $v, w \in \mathscr{S}$. Furthermore if some index in $s$ is 1, then (49) implies that

$$\|\Pi_M(s)v - \Pi_M(s)w\|_1 \leqslant r'\|v - w\|_1. \tag{50}$$

For fixed $v, w \in \mathscr{S}$, let s be the $(N - 1)$-string described in Assumption 3.4(B), and let $s'$ be the $N$-string obtained by adding the last index $i_N = 1$ to s, so that all flows are dropped at the last congestion event in this sequence. Then Assumption 3.4 implies that $p^*(s'; w) \geqslant qq'$. Combining this with (50) we

deduce that

$$\sum_s p^*(s'; w)\|\Pi_N(s)v - \Pi_N(s)w\|_1$$

$$= p^*(s'; w)\|\Pi_N(s')v - \Pi_N(s')w\|_1$$
$$+ \sum_{s \neq s'} p^*(s; w)\|\Pi_N(s)v - \Pi_N(s)w\|_1$$
$$\leqslant (r' p^*(s'; w) + 1 - p^*(s'; w))\|v - w\|_1$$
$$\leqslant (1 - (1 - r')qq')\|v - w\|_1 \tag{51}$$

and this establishes (46). To prove (47), let $G = \sup_{v,w \in \mathcal{S}} \|v - w\|_1$, and take $M'$ large enough so that

$$(r')^{M'} G < \varepsilon. \tag{52}$$

We lower bound the left side of (47) by the single term $p^*(s'; v)p^*(t'; w)$. Here $s'$ and $t'$ are $(N - 1 + M')$-strings. The first $N - 1$ indices in $s'$ are those described in Assumption 3.4(B), where the starting point is $v$. The remaining $M'$ indices are all 1, corresponding to repeated applications of the matrix $A_1$. It follows from Assumption 3.4 that $p^*(s'; v) \geqslant q'q^{M'}$. The string $t'$ is defined similarly with starting point $w$. Then taking $M = M' + N - 1$ gives (47) with $\delta = (q'q^{M'})^2$. $\quad\square$

We now apply Theorem 2.1 of Barnsley et al. (1988), which states that under the conditions of Lemma 4.3 there is an attractive, unique invariant probability measure for the Markov chain. In their proof of this result, Barnsley et al. assume average contractivity at each step of the process, that is they assume $N = 1$. However, as they remark the proof extends to the case where there is average contractivity over some fixed number of steps $N$ (independent of the initial points) as in our Lemma 4.3. In fact as noted before the proof in our case is simpler, as we have a compact state space and we assume uniform Lipschitz regularity for the drop probability functions. This establishes (i) and (ii). Property (iii) follows by Elton's (1990) result, which states that time averages converge almost surely to ensemble averages for this model.

## References

Barnsley, M., Demko, S., Elton, J., & Geronimo, J. (1988). Invariant measures for Markov processes arising from iterated function systems with place dependent probabilities. *Annales de l'Institut Henri Poincaire Probability and Statistics*, *24*, 367–394.

Baccelli, F., & Hong, D. (2002). AIMD, fairness and fractal scaling of TCP traffic. In *Proceedings of IEEE INFOCOM 2002*, June 2002. (pp. 229–238). New York, NY, USA.

Berman, A., & Plemmons, R. (1979). *Nonnegative matrices in the mathematical sciences*. Philadelphia, PA: SIAM.

Chiu, D., & Jain, R. (1989). Analysis of the increase/decrease algorithms for congestion avoidance in computer networks. *Journal of Computer Networks*, *17*, 1–14.

Elton, J. (1990). A multiplicative ergodic theorem for Lipschitz maps. *Stochastic Processes Applications*, *231*(2), 39–47.

Hespanha, J. (2004). Stochastic hybrid systems: Application to communication networks. In R. Alur, & G.J. Pappas (Eds.), *Hybrid systems*, *computation & control*, *Proceedings of HSCC 2004*, *Lecture notes in computer science*, March 2004. (Vol. 2993, pp. 387–401). Heidelberg: Springer.

Isaac, R. (1962). Markov processes and unique probability measures. *Pacific Journal of Mathematics*, *12*, 273–286.

Leizarowitz, A., Stanojevic, R., & Shorten, R. (2006). Tools for the analysis and design of communication networks with Markovian dynamics. *IEE Proceedings on control theory*, in press.

Padhye, J., Firoiu, V., Towsley, D. F., & Kurose, J. F. (2000). Modeling TCP Reno performance: A simple model and its empirical validation. *IEEE/ACM transactions on networking*, *8*(2), 133–145.

Shorten, R. N., Leith, D. J., Foy, J., & Kilduff, R. (2005). Analysis and design of AIMD congestion control algorithms in communication networks. *Automatica*, *41*, 725–730.

Shorten, R., Wirth, F., & Leith, D. (2006). A positive systems model of TCP-like congestion control: Asymptotic analysis. *IEEE/ACM transactions on networking*, *14*, 616–629.

Srikant, R. (2004). *Internet congestion control, Control theory Vol. 14*, Boston, MA: Birkhäuser Boston Inc.

Stenflo, O. (2002). Uniqueness of invariant measures for place dependent iterations of functions. *IMA Journal of Applied Mathematics*, *132*, 13–32.

Wirth, F., Stanojevic, R., Shorten, R., & Leith, D. (2006). Stochastic equilibria of AIMD communication networks. *SIAM Journal on Matrix Analysis and Applications*, in press.

**Robert Shorten** graduated from the University College Dublin (UCD) in 1990 with a First Class Honours B.E. degree in Electronic Engineering. From 1993 to 1996 Dr. Shorten was the holder of a Marie Curie Fellowship to conduct research at the Daimler-Benz Research Institute for Information Technology in Berlin. In 1997 Dr. Shorten was awarded a European Presidential fellowship to return to Ireland. Prof. Shorten is a co-founder and a Senior Researcher of the Hamilton Institute, NUI Maynooth and is an Editor of the IEE Processings on Control Theory. His research interests include stability theory, linear algebra, and network congestion control.

**Christopher King** Christopher King was born in Dublin, Ireland. He received the B.A. degrees in Mathematics and Physics at Trinity College Dublin in 1980, and the Ph.D. degree at Harvard University, Cambridge, MA in 1984. He is currently Professor of Mathematics at Northeastern University in Boston, MA. His research interests include mathematical physics, quantum information theory and network analysis, and he is a member of the editorial board of the Journal of Mathematical Physics. His research in quantum information theory is supported by the National Science Foundation.

**Fabian Wirth** received his Diploma, Dr. rer. nat. and venia legendi in mathematics from the University of Bremen, Germany, where he was with the Institute for Dynamical Systems and the Centre for Technomathematics. He is senior researcher at the Hamilton Institute, where he works on the dynamics of communication networks. His interests include stability theory of dynamical systems and robust stability.

**Douglas Leith** graduated from the University of Glasgow in 1986 with a first class B.Sc. (Eng) degree in Electronics and Electrical Engineering and Computer Science and was awarded his Ph.D. also from the University of Glasgow, in 1989. Following the award of a Royal Society personal research fellowship to study nonlinear control, in 2001 Prof. Leith joined the National University of Ireland Maynooth as Director of the Hamilton Institute (www.hamilton.ie). Prof Leith's current research interests include internet congestion control and dynamics, resource allocation in wireless networks and nonlinear time series analysis.