

# A network flow model for mixtures of file transfers and streaming traffic

Peter Key, Laurent Massoulié

Microsoft Research, 7 J J Thompson Avenue, Cambridge CB3 0FB, UK

Alan Bain, Frank Kelly

Statistical Laboratory, University of Cambridge, Cambridge CB3 0WB, UK

Technical Report

MSR-TR-2003-37

Microsoft Research  
Microsoft Corporation  
One Microsoft Way  
Redmond, WA 98052

<http://www.research.microsoft.com>

## Abstract

Roberts, Massoulié and co-authors have introduced and studied a flow level model of Internet congestion control, that represents the randomly varying number of flows present in a network where bandwidth is dynamically shared between elastic file transfers. In this paper we consider a generalization of the model to include streaming traffic as well as file transfers, under a fairness assumption that includes TCP-friendliness as a special case. We establish stability, under conditions, for a fluid model of the system. We also assess the impact of each traffic type on the other: file transfers are seen by streaming traffic as reducing the available capacity, whereas for file transfers the presence of streaming traffic amounts to replacing sharp capacity constraints by relaxed constraints. The integration of streaming traffic and file transfers has a stabilizing effect on the variability of the number of flows present in the system.

## 1 INTRODUCTION

The current Internet is dominated by flows which use TCP. The percentage of TCP traffic is variable, and may depend on time of day and the particular route chosen; however typical measurements on a backbone [11] show that upwards of 70% of flows use TCP, rising to over 90% by volume, with UDP the main alternative protocol (up to 20% of packets, or 10% of bytes). Prevailing applications can change rapidly: whereas Web traffic used to be the dominant application type for TCP traffic, at the time of writing file-sharing applications can dominate. The current volume of streaming traffic carried by UDP is small (less than 10%), but the rapid increase in peer-to-peer traffic illustrates how quickly the status-quo can change, and we would like to predict behaviour in different scenarios.

How TCP and UDP should co-exist is a vexed question, and many regard UDP-related traffic as inherently problematic. Some authors have proposed that streaming traffic should be TCP-friendly, so that it can share network resources fairly with the dominant form of existing traffic [9]. Applications that use UDP often need some form of quality of service to function adequately, which has led some researchers to consider distributed or end-point admission control [12, 4, 1].

We are motivated by the need to model such situations, which requires modelling the heterogeneous traffic streams, with their different characteristics. We consider two types of traffic, which we label ‘file transfers’ and ‘streaming’ traffic. A flow carrying a file transfer must transfer a given volume: the volume may be random, but is independent of the level of congestion experienced. An admitted flow carrying streaming traffic remains present for a holding time: the holding time may be random, but is independent of the level of congestion experienced.

The analysis of streaming traffic on its own gives rise to a product-form solution under certain reasonable assumptions, a form which is preserved under certain types of call admission control [12]. Moreover the limiting behaviour as the size of the system grows leads naturally to a non-degenerate limit for the (scaled) number of connections. In contrast, a similar scaling applied to just file transfer traffic leads to a distribution that is either unstable or has mean zero; it has been suggested [6] that such a model is flawed, lacking any self-limiting behaviour. We shall see that this criticism is avoided when the two types of traffic are mixed, and that the presence of even a small amount of streaming traffic has a stabilising effect.

The organization of this paper is as follows. In Section 2 we describe the bandwidth sharing policy we assume, a generalized form of TCP-friendliness. In Section 3 we describe the flow level model, a generalization of the model introduced by Massoulié and Roberts [14]. In Section 4 we establish stability, under conditions, for a fluid model of the system, through the construction of an appropriate Lyapunov function. In Section 5 we consider extensions, including admission control. In Section 6 we discuss simulations of the flow level model for a star network, and explore the impact of streaming traffic on the variability of flow bandwidth. We conclude in Section 7.

## 2 FAIRNESS ASSUMPTIONS

Consider a network with resources labelled by  $j \in J$ . Let a route  $r$  identify a non-empty subset of  $J$  (interpreted as the set of resources used by a flow on route  $r$ ). Write  $R$  for the set of possible routes. Set  $A_{jr} = 1$  if resource  $j$  lies on route  $r$  (i.e.  $j \in r$ ), and set  $A_{jr} = 0$  otherwise. We assume positive finite capacities ( $C_j, j \in J$ ).

Let  $N_r$  be the number of flows on route  $r$ . Given a fixed parameter  $\alpha \in (0, \infty)$  and strictly positive weights ( $w_r, r \in R$ ), we suppose that the bandwidth allocation to each of the  $N_r$  flows on route  $r$  is  $x_r$ , where  $x = (x_r, r \in R)$  is a solution to the following optimization problem:

$$\text{maximize} \quad \sum_{r \in R} w_r N_r \frac{x_r^{1-\alpha}}{1-\alpha} \quad (1)$$

$$\text{subject to} \quad \sum_r A_{jr} N_r x_r \leq C_j, \quad j \in J \quad (2)$$

$$\text{over} \quad x_r \geq 0, \quad r \in R. \quad (3)$$

Call the resulting allocation a weighted  $\alpha$ -fair allocation [16].

The form of a solution to the problem (1–3) can be given in terms of Lagrange multipliers ( $p_j, j \in J$ ), one for each of the capacity constraints (2), as

$$x_r = \left( \frac{w_r}{\sum_j p_j A_{jr}} \right)^{1/\alpha}, \quad (4)$$

where

$$p_j \geq 0, \quad p_j \left( C_j - \sum_r A_{jr} N_r x_r \right) = 0 \quad j \in J. \quad (5)$$

The strict concavity of the objective function (1) as a function of  $(x_r, r : N_r > 0)$  ensures that the component  $x_r$  is unique if  $N_r$  is positive. When  $w_r = 1, r \in R$ , the cases  $\alpha \rightarrow 0$ ,  $\alpha \rightarrow 1$  and  $\alpha \rightarrow \infty$  correspond respectively to an allocation which achieves maximum throughput, is *proportionally fair* or is *max-min fair* [3, 16]. Weighted  $\alpha$ -fair allocations provide a tractable theoretical abstraction of decentralized packet-based congestion control algorithms such as TCP.

If  $\alpha = 2$  and  $w_r$  is the reciprocal of the square of the round trip time on route  $r$ , then the formula (4) is a version of the *inverse square root law* familiar from studies of the throughput of TCP connections [8, 15, 17]. A flow carrying streaming traffic is termed *TCP-friendly* if, *inter alia*, it adapts its rate to correspond with the steady-state rate of a TCP connection, usually characterized in terms of a version of the inverse square root law [9].

The relations (2–5), and more refined versions of these relations, can be solved to give predictions of throughput, given the numbers of flows  $N$  present [5, 10, 19]. Given  $N$ , network performance along different routes can be predicted. But what determines the behaviour of  $N$ ? One aim of this paper is to better understand how the behaviour of  $N$  is influenced by the mix of traffic types present.

## 3 FLOW LEVEL STOCHASTIC MODEL

We now describe our model of how flows arrive and depart. Our aim is to generalize the stochastic model for file transfers introduced in [14] to include streaming flows.

Let  $N_r$  be the number of document transfers on route  $r$ , and let  $M_r$  be the number of streaming flows on route  $r$ . Define the indicator function  $I[r = s] = 1$  if  $r = s$ ,  $I[r = s] = 0$  otherwise. Let  $T_s N = (N_r + I[r =$

$s], r \in R)$ , with inverse  $T_s^{-1}N = (N_r - I[r = s], r \in R)$ . We suppose that  $(N, M) = (N_r, r \in R; M_r, r \in R)$  is a Markov process, with state space  $\mathbb{Z}_+^J \times \mathbb{Z}_+^J$  and non-trivial transition rates

$$q((N, M), (T_r N, M)) = \nu_r, \quad q((N, M), (T_r^{-1} N, M)) = \mu_r N_r x_r(N + M), \quad r \in R$$

$$q((N, M), (N, T_r M)) = \kappa_r, \quad q((N, M), (N, T_r^{-1} M)) = M_r \eta_r, \quad r \in R$$

for  $(N, M) \in \mathbb{Z}_+^J \times \mathbb{Z}_+^J$ , where  $x(N)$  is a solution to the optimization problem (1–3). This corresponds to a model where new file transfers arrive on route  $r$  as a Poisson process of rate  $\nu_r$ , new streaming flows arrive on route  $r$  as a Poisson process of rate  $\kappa_r$ , and  $x_r(N + M)$  is the bandwidth allocated to each flow on route  $r$ , whether it is a file transfer or streaming flow. A file transfer on route  $r$  transfers a file whose size is exponentially distributed with parameter  $\mu_r$ , and a streaming flow on route  $r$  has an exponentially distributed holding time with parameter  $\eta_r$ .

If  $\kappa_r = 0, r \in R$ , then this model reduces to the model introduced by Roberts and Massoulié [14], in which there are no streaming flows, only file transfers. For this case, De Veciana, Lee and Konstantopoulos [7] and Bonald and Massoulié [3] have shown that a sufficient condition for the Markov chain  $(N(t), t \geq 0)$  to be positive recurrent is that

$$\sum_r A_{jr} \rho_r < C_j, \quad j \in J, \quad (6)$$

where  $\rho_r = \nu_r / \mu_r$ ; this condition is also necessary [13]. The condition is natural:  $\rho_r$  is the load on route  $r$ , and we can identify the ratio of the two sides of the inequality (6) as the *traffic intensity* at resource  $j$ . Kelly and Williams [13] have explored the behaviour of a fluid model for this case in heavy traffic, when the inequalities (6) are close to being tight, which is a key step towards proving state space collapse. The papers [3, 7, 13] all make use of a fluid model of the Markov process, an approach which we shall use for our analysis of the extended model.

We shall henceforth assume that  $\kappa_r > 0, r \in R$ , and that condition (6) is satisfied. Define the *reduced capacities*

$$\tilde{C}_j = C_j - \sum_r A_{jr} \rho_r, \quad j \in J. \quad (7)$$

Thus the reduced capacity  $C_j$  on resource  $j$  is just the amount by which inequality (6) fails to be tight. The reduced capacities will determine the capacity available to streaming flows in a sense that will be made precise in the next section.

## 4 STABILITY OF FLUID MODELS

Next we describe a fluid model, which can be thought of as a formal law of large numbers approximation under the scaling

$$(n, m)(t) = \left( \frac{N_c(t)}{c}, \frac{M_c(t)}{c} \right) \quad c \rightarrow \infty,$$

where  $(N_c(t), M_c(t))$  is the model of the previous Section but with  $C_j, j \in J$ , and  $\nu_r, \kappa_r, r \in R$ , replaced by  $cC_j, j \in J$ , and  $c\nu_r, c\kappa_r, r \in R$ , respectively. The fluid model is an approximation appropriate for the case where  $C_j, j \in J$ , and  $\nu_r, \kappa_r, r \in R$ , are all large, an important case in applications.

The fluid model for the Markov process of the last Section takes the form

$$\frac{d}{dt} n_r(t) = \nu_r - \mu_r n_r(t) x_r(n(t) + m(t)), \quad r \in R \quad (8)$$

$$\frac{d}{dt} m_r(t) = \kappa_r - \eta_r m_r(t), \quad r \in R. \quad (9)$$

Note that our assumption that  $\kappa_r > 0, r \in R$ , implies that  $m_r(t) > 0, r \in R, t > 0$ .

**Proposition 1.** *Provided the condition (6) is satisfied, the differential equations (8,9) have a unique invariant point,  $(\hat{n}_r, \hat{m}_r)$ . It takes the form  $\hat{m}_r = \kappa_r/\eta_r$  and*

$$\hat{n}_r = \frac{\nu_r}{\mu_r} \left( \frac{\sum_{j \in J} p_j A_{jr}}{w_r} \right)^{1/\alpha} \quad r \in R, \quad (10)$$

for some  $p \in \mathbb{R}_+^J$ . At the invariant point the bandwidth allocation to each flow on route  $r$  is

$$x_r = \left( \frac{w_r}{\sum_j p_j A_{jr}} \right)^{1/\alpha}. \quad (11)$$

The pair  $(x, p)$  forms a solution of equation (11) and the conditions

$$p_j \geq 0, \quad p_j \left( \tilde{C}_j - \sum_r A_{jr} \hat{m}_r x_r \right) = 0 \quad j \in J, \quad (12)$$

and together these relations determine  $x$  uniquely.

*Proof.* At an invariant point  $m_r(t) = \hat{m}_r$ , from equation (9). Further,

$$\hat{n}_r x_r (\hat{n} + \hat{m}) = \rho_r, \quad (13)$$

from equation (8). Now at any time  $t$ ,

$$x_r(n(t) + m(t)) = \left( \frac{w_r}{\sum_j p_j(t) A_{jr}} \right)^{1/\alpha}$$

where

$$p_j(t) \geq 0, \quad p_j(t) \left( C_j - \sum_r A_{jr} (n_r(t) + m_r(t)) x_r(n(t) + m(t)) \right) = 0 \quad j \in J,$$

from the characterization of  $x$  as a solution to an optimization problem of the form (1–3). Thus, at an invariant point,

$$p_j \geq 0, \quad p_j \left( \tilde{C}_j - \sum_r A_{jr} \hat{m}_r \left( \frac{w_r}{\sum_j p_j A_{jr}} \right)^{1/\alpha} \right) = 0 \quad j \in J,$$

using equation (13) and the definition (7). Thus  $x$ , given by (11), is the unique optimum to a problem of the form (1–3), with  $C$  replaced by  $\tilde{C}$  and  $N$  replaced by  $\hat{m}$ .  $\square$

Equation (10) describes the vector  $\hat{n}$ , of dimension  $|R|$ , in terms of  $p$ , a vector which may have a much smaller dimension,  $|J|$ , a phenomenon first noted in the balanced fluid model of [13].

The invariant point can be interpreted as follows. File transfers place an irreducible load  $\sum_r A_{jr} \rho_r$  on resource  $j$  for each  $j \in J$ . The reduced capacities  $(\tilde{C}_j, j \in J)$  that remain after this load is satisfied are available to be shared amongst streaming traffic, and determine the bandwidth allocation to flows on route  $r$  for both types of traffic.

When  $\kappa_r = 0, r \in R$ , the unique invariant point of the fluid model is  $\hat{n} = 0$  [7, 3]. It is notable that the inclusion of streaming traffic within the fluid model forces the components of  $\hat{n}$  to be positive.

We now discuss convergence to the equilibrium point of the above dynamics. In order to do so, it is convenient to introduce a modification for the dynamics of file transfers. This is naturally described in terms of the quantities

$\lambda_r$ , which represent the total capacity allocated to type  $r$  file transfers, and thus with the previous notation,  $\lambda_r = n_r x_r$ . Let the function  $\psi(\lambda)$  be a penalty function. Then the modified dynamics are as follows:

$$\frac{d}{dt}n_r(t) = \nu_r - \mu_r \lambda_r(n(t)), \quad r \in R, \quad (14)$$

where the vector  $\lambda$  of service rates  $\lambda_r$  is defined as the solution to the optimisation problem

$$\text{maximize} \quad \phi(\lambda) := \sum_{r \in R} w_r n_r^\alpha \frac{\lambda_r^{1-\alpha}}{1-\alpha} + \psi(\lambda) \quad (15)$$

$$\text{subject to} \quad \sum_r A_{jr} \lambda_r \leq C_j, \quad j \in J \quad (16)$$

$$\text{over} \quad \lambda_r \geq 0, \quad r \in R. \quad (17)$$

In the case where  $\psi$  is identically zero, this reduces to the previous dynamics for the file transfers in the absence of streaming traffic. The function  $\psi$  is assumed to be concave and strictly monotonic decreasing in each coordinate on the domain of the optimisation problem. This latter condition implies that the rate  $\lambda_r$  goes to zero as  $n_r$  goes to zero, and hence the trajectories  $n_r$  stay away from the boundary of the orthant  $\mathbb{R}_+^R$ . Let us prove stability of the above dynamics.

**Theorem 2.** *Under the stability conditions (6), the function  $L(n)$  defined by*

$$L(n) = \sum_r \frac{1}{\mu_r} \left\{ w_r \frac{n_r^{1+\alpha}}{(1+\alpha)\rho_r^\alpha} + n_r \psi'_r(\rho) \right\}, \quad (18)$$

where  $\psi'_r(\rho)$  stands for the  $r$ -th partial derivative  $\frac{\partial \psi}{\partial \lambda_r}$  evaluated at the vector of loads  $\rho_r$ , is a Lyapunov function for the dynamics (14–17). Hence these dynamics converge to the unique minimiser of  $L$  on the orthant  $\mathbb{R}_+^R$ , that is

$$\hat{n}_r = \rho_r \left( \frac{-\psi'_r(\rho)}{w_r} \right)^{1/\alpha}. \quad (19)$$

*Proof.* Under the condition (6), the vector  $\rho = (\rho_r, r \in R)$  lies in the interior of the domain (16–17) of the optimisation problem defining the vector  $\lambda$ . Since the function  $\phi$  is strictly concave on this domain\*, it holds that

$$\sum_r \phi'_r(\rho)(\rho_r - \lambda_r) \leq 0,$$

and this inequality is strict unless  $\lambda = \rho$ . The left-hand side also reads

$$\sum_r \left\{ w_r \left( \frac{n_r}{\rho_r} \right)^\alpha + \psi'_r(\rho) \right\} (\rho_r - \lambda_r),$$

and is thus equal to

$$\sum_r \frac{\partial L}{\partial n_r}(n(t)) \frac{d}{dt}n_r(t) = \frac{d}{dt}L(n(t)).$$

Thus the value of  $L(n)$  decreases strictly along the trajectories of the system, except at the equilibrium point specified by (19), which is the only point for which the corresponding rate vector  $\lambda$  equals the load vector  $\rho$ .  $\square$

---

\*Strict concavity of  $\phi$  follows from concavity of the two terms in its definition (15) and strict concavity of the first term in (15).

**Remark 3.** If the concave function  $\psi$  fails to be differentiable at  $\rho$ , by adapting the above proof it can be shown that the dynamics (14–17) converge to the set of points  $\hat{n}$  satisfying (19), where the vector  $(-\psi'_r(\rho), r \in R)$  spans the set of sub-gradients of the convex function  $-\psi$  at  $\rho$ . We refer the reader to [18], p.214 for a definition and basic properties of sub-gradients of convex functions.

We now apply this result to establish stability of the dynamics (8–9).

**Corollary 4.** Under the stability condition (6), the dynamics (8–9) are asymptotically stable.

*Proof.* We shall only treat the special case where the  $m_r$  have already converged to their equilibrium values,  $\hat{m}_r$ . As the convergence of  $m(t)$  to  $\hat{m}$  does not depend on the evolution of  $n(t)$ , the general case can be deduced by continuity arguments. We now show that the  $n_r$  evolve according to (14–17) for some suitable choice of a penalty function  $\psi$ . Indeed, (14) holds, with the service rates  $\lambda_r$  solving

$$\begin{aligned} & \text{maximize} && \phi(\lambda, \gamma) := \sum_{r \in R} w_r \left\{ n_r^\alpha \frac{\lambda_r^{1-\alpha}}{1-\alpha} + \hat{m}_r^\alpha \frac{\gamma_r^{1-\alpha}}{1-\alpha} \right\} \\ & \text{subject to} && \sum_r A_{jr}(\lambda_r + \gamma_r) \leq C_j, \quad j \in J \\ & \text{over} && \lambda_r, \gamma_r \geq 0, \quad r \in R. \end{aligned}$$

Performing the optimisation over the  $\gamma_r$  first, this is of the form (15–17), with

$$\begin{aligned} \psi(\lambda) &:= \sup \left\{ \sum_r w_r \hat{m}_r^\alpha \frac{\gamma_r^{1-\alpha}}{1-\alpha} \right\}, \\ \text{over } \gamma \in S(\lambda) &:= \left\{ \gamma \in \mathbb{R}_+^R, \sum_r A_{jr} \gamma_r \leq C_j - \sum_r A_{jr} \lambda_r, \quad j \in J \right\}. \end{aligned} \tag{20}$$

It is readily seen that  $\psi$  is decreasing in each coordinate: given  $\lambda, \lambda'$ , such that  $\lambda'_r \leq \lambda_r$  for all  $r$ , the inequality being strict for some  $r$ , any vector  $\gamma$  in  $S(\lambda)$  is such that  $\gamma' := (\gamma_r + \lambda_r - \lambda'_r)$  is in  $S(\lambda')$ , so that  $\psi(\lambda) < \psi(\lambda')$ . Concavity of  $\psi$  also holds: given  $\lambda, \lambda'$  and  $\epsilon$  in  $[0, 1]$ , denote by  $\gamma$  and  $\gamma'$  the maximising vectors in the definition of  $\psi(\lambda), \psi(\lambda')$  respectively. Then  $\epsilon\gamma + (1-\epsilon)\gamma'$  lies in  $S(\epsilon\lambda + (1-\epsilon)\lambda')$ , and hence

$$\psi(\epsilon\lambda + (1-\epsilon)\lambda') \geq \sum_r w_r \hat{m}_r^\alpha \frac{(\epsilon\gamma_r + (1-\epsilon)\gamma'_r)^{1-\alpha}}{1-\alpha} \geq \epsilon\psi(\lambda) + (1-\epsilon)\psi(\lambda'),$$

where concavity of the function maximised in the definition of  $\psi$  gives the second inequality.  $\square$

**Remark 5.** Under the particular choice (20) of penalty function, and comparing equations (10) and (19), we deduce that  $\sum_{j \in J} p_j A_{jr} = -\psi'_r(\rho)$ . Notice the identification between the sensitivity of the penalty function  $\psi$  with respect to the load  $\rho_r$  and the sum of the Lagrange multipliers along route  $r$ .

## 5 EXTENSIONS: PACKET MODELS AND ADMISSION CONTROL

### 5.1 Constraint relaxation

The formulation (14–17) is also useful to model situations where the hard capacity constraints described by the intersection of half-spaces (2) are relaxed. If the optimization problem (1–3) is replaced by

$$\begin{aligned} & \text{maximize} && \sum_r w_r N_r \frac{x_r^{1-\alpha}}{1-\alpha} - \sum_j C_j \left( \sum_r A_{jr} N_r x_r \right) \\ & \text{over} && x \geq 0, \end{aligned}$$

where  $C_j(\cdot), j \in J$ , are convex, strictly increasing, differentiable functions, then an optimum is again given by equation (4), but where now  $p_j, j \in J$ , satisfy

$$p_j = C'_j \left( \sum_r A_{jr} N_r x_r \right).$$

This formulation arises naturally from packet level models, with  $x_r$  the mean rate of a stochastic packet generation process. For example, if the resources  $j$  correspond to output ports of routers, then there is a limited amount of buffering available, and packets will be dropped if the capacity is exceeded, or more generally marked according to some active queue management technique. We may interpret  $p_j(y_j)$  as the probability of dropping (or marking) a packet at resource  $j$  when the load on the resource is  $y_j$ .

Stability of the corresponding fluid model can be deduced from the formulation (14–17), by setting

$$\psi(\lambda) = - \sum_j C_j \left( \sum_r A_{jr} \lambda_r \right).$$

## 5.2 Admission controlled traffic

Streaming may need some minimal non-zero rate for the application to function adequately. For example in the case of streaming multimedia, even with adaptive codecs, some minimal transmission rate is often required for acceptable performance. Suppose that type  $r$  streaming traffic only enters if  $x_r \geq x_r^{min}$ : then in both the flow level stochastic model and in the fluid limit,  $\kappa_r$  is replaced by  $\kappa_r I[x_r \geq x_r^{min}]$ . At an invariant point, either  $m_r > 0$  and  $x_r \geq x_r^{min}$  or  $m_r = 0$ . The condition  $x_r \geq x_r^{min}$  is equivalent to

$$\sum_j p_j A_{jr} \leq \frac{w_r}{(x_r^{min})^\alpha} \quad r \in R. \quad (21)$$

Thus an invariant point  $(\hat{n}_r, \hat{m}_r)$  is described in terms of a vector  $p$  of dimension  $|J|$  which lies in the polyhedral region defined by the intersection of the positive orthant with the  $|R|$  half-spaces (21). If the parameters  $p_j, j \in J$ , satisfy the linear constraints (21) with strict inequality, then the fluid model predicts there will be no call admission blocking.

## 6 EXAMPLE: A STAR NETWORK

As a concrete example, consider a star network of 10 links connected to a core. This example is motivated by the current Internet, where the back-bone is relatively uncongested, and congestion occurs mainly on the access links. Flows use two links, with traffic spread randomly across links.

For the example,  $J = \{1, 2, \dots, 10\}$ ,  $R = \{(i, j) : i < j, i, j \in J\}$ . The capacity of each link  $C_j$  was chosen equivalent to a T3 link (45 Mbit/s), for  $j \in J$ . The mean holding time of streaming traffic ( $1/\eta_r, r \in R$ ) was taken to be 200 seconds, corresponding to voice traffic, with the mean file size ( $1/\mu_r, r \in R$ ) taken to be 600kB. The arrival rates for the two types of traffic ( $\nu_r$ , and  $\kappa_r$ ) were chosen to be identical, giving a file-transfer traffic intensity of 0.5 on each link, and such that in equilibrium each flow has rate 25kbit/s ( $x_r$ ). Under this regime the equilibrium number of flows of each type is 100 ( $\hat{m}_r = \hat{n}_r = 100$ ) per route  $r$ , giving 900 flows of each type on each link  $j$ .

Figure 1(a) shows the evolution of the number of each type  $\sum_r A_{jr} n_r$  and  $\sum_r A_{jr} m_r$  on a typical *link*, obtained by simulation of the Markov chain of Section 3. Note that the two curves look very similar and have a

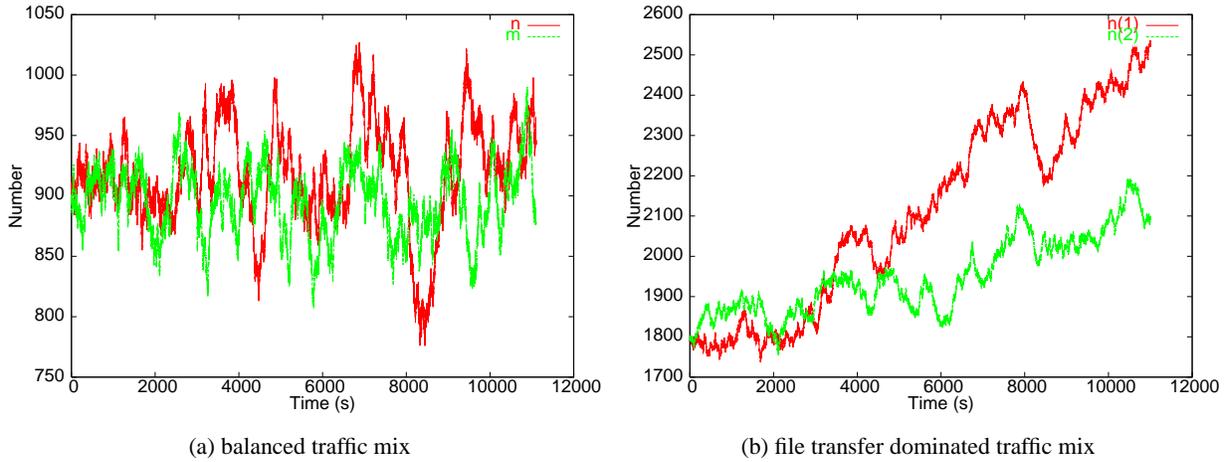


Figure 1: Impact of streaming traffic on file transfers. The substantial amount of streaming traffic present in mix (a) relative to mix (b) has a stabilizing effect on the number of flows in progress.

mean of 900. The number of streaming flows in progress has a standard deviation of 30, while the number of file transfers has a slightly higher standard deviation of just over 40.

We now alter the offered load of each type of traffic, to keep the nominal quality ( $x_r$ ) seen by the flows fixed while significantly altering the proportions of the two types of traffic. We make the file-transfer traffic intensity 0.995 on each link, with a very small amount of streaming traffic. The load was such that the equilibrium of the fluid model has  $\hat{n}_r = 199, \hat{m}_r = 1$ . (With so little streaming traffic we do not expect our fluid model to be a good approximation; as the amount of streaming traffic decreases to zero, we expect the behaviour of the system to be better described by the Brownian model of [13].) In Figure 1(b) we plot the behaviour of  $\sum_r A_{jr} n_r$  on two typical links: observe the different vertical scale in this figure, and the marked variability of the number of flows in progress. Comparing the two figures, we see that the substantial proportion of streaming traffic present, in Figure 1(a), has the effect of reducing the variability of the number of flows in progress, and hence the variability of the bandwidth received by flows.

## 7 CONCLUSION

We have studied a flow level model of Internet congestion control, that represents the randomly varying number of flows present in a network. Bandwidth was assumed to be dynamically shared between file transfers and streaming traffic, according to a fairness criterion that includes TCP friendliness as a special case. Through the construction of an appropriate Lyapunov function we have established stability, under conditions, for a fluid model of the system. The presence of fair-sharing streaming traffic results in a non-degenerate fluid model. Analysis of the model suggests that file transfers are seen by streaming traffic as reducing the available capacity, whereas for file transfers the presence of streaming traffic amounts to replacing sharp capacity constraints by relaxed constraints. Simulations show that the integration of streaming traffic and file transfers has a stabilizing effect on the variability of the number of flows present in the system.

## References

- [1] A. Bain and P. B. Key (2001, Modelling the performance of distributed admission control for adaptive applications, *Performance Evaluation Review*, December 2001.

- [2] S. Ben Fredj, T. Bonald, A. Proutiere, G. Regnie, J. Roberts (2001), Statistical bandwidth sharing: a study of congestion at flow level. In *Proceedings of SIGCOMM 2001*.
- [3] T. Bonald and L. Massoulié (2001), Impact of fairness on Internet performance. In *Proceedings of ACM SIGMETRICS 2001*.
- [4] L. Breslau, E. W. Knightly, S. Shenker, I. Stoica, and H. Zhang (2000), Endpoint admission control: Architectural issues and performance. In *Proceedings of SIGCOMM 2000*, 57–69.
- [5] T. Bu and D. Towsley (2001), Fixed Point Approximation for TCP behavior in an AQM Network. In *Proceedings of ACM SIGMETRICS 2001*.
- [6] C. A. Courcoubetis, A. Dimakis, and M. I. Reiman (2001), Providing bandwidth guarantees over a best-effort network: call admission and pricing. *IEEE INFOCOM*, 459–467.
- [7] G. de Veciana, T.-J. Lee and T. Konstantopoulos (2001), Stability and performance analysis of networks supporting elastic services, *IEEE/ACM Trans. on Networking* **9**, 2–14.
- [8] S. Floyd and K. Fall (1999), Promoting the use of end-to-end congestion control in the Internet. *IEEE/ACM Transactions on Networking* **7**, 458–472.
- [9] Floyd, S., M. Handley, J. Padhye and J. Widmer (2000), Equation-based congestion control for unicast applications. In *Proc. ACM SIGCOMM 2000*, pages 43-54, Stockholm.
- [10] R.J. Gibbens, S.K. Sargood, C. Van Eijl, F.P. Kelly, H. Azmoodeh, R.N. Macfadyen and N.W. Macfadyen (2000), Fixed-point models for the end-to-end performance analysis of IP networks. 13th ITC Specialist Seminar: IP Traffic Measurement, Modeling and Management, Sept 2000, Monterey, California.
- [11] IP monitoring project, Sprint labs. <http://ipmon.sprintlabs.com>
- [12] F.P. Kelly, P.B. Key, and S. Zachary (2000), Distributed admission control. *IEEE Journal on Selected Areas in Communications*, **18**, 2617–2628.
- [13] F.P. Kelly and R. J. Williams (2003), Fluid model for a network operating under a fair bandwidth-sharing policy. <http://www.math.ucsd.edu/~williams>
- [14] L. Massoulié and J. Roberts (2000), Bandwidth sharing and admission control for elastic traffic. *Telecommunication Systems* **15**, 185–201.
- [15] M. Mathis, J. Semke, J. Mahdavi, and T. Ott (1997), The macroscopic behaviour of the TCP congestion avoidance algorithm. *Computer Communication Review* **27**, 67–82.
- [16] J. Mo and J. Walrand (2000), Fair end-to-end window-based congestion control. *IEEE/ACM Transactions on Networking* **8**, 556–567.
- [17] J. Padhye, V. Firoiu, D. Towsley and J. Kurose (2000), Modeling TCP Reno performance: a simple model and its empirical validation. *IEEE/ACM Transactions on Networking* **8**, 133–145.
- [18] T. Rockafellar (1970), *Convex Analysis*. Princeton University Press.
- [19] M. Roughan, A. Erramilli and D. Veitch (2001), Network performance for TCP Networks, Part I: persistent sources. In *Proceedings of ITC'17 Brasil*, September 2001.