# On the Stability Region of Congestion Control

Xiaojun Lin and Ness B. Shroff

School of Electrical and Computer Engineering

Purdue University, West Lafayette, IN 47906

{linx,shroff}@ecn.purdue.edu

**Abstract**

It is well known that congestion control can be viewed as a distributed iterative algorithm solving a global optimization problem that maximizes the total system utility. In this paper, we study the stability region of a network employing congestion control algorithms derived from such an optimization framework. Previous work in the literature typically adopts a *time-scale separation assumption*, which assumes that, whenever the number of users in the system changes, the data rates of the users are adjusted *instantaneously* to the optimal rate allocation computed by the global optimization problem. Under this assumption, it has been shown that such rate allocation policies can achieve the largest possible stability region. However, this time-scale separation assumption, although technically convenient, rarely holds in practice. In this paper, we remove this time-scale separation assumption and show that the largest possible stability region can still be achieved by a large class of congestion control algorithms derived from the optimization framework. Our result provides new insights on the performance implication of congestion control and on the choices of the parameters of the congestion controller.

# 1 Introduction

Congestion control (or rate control) is a key functionality in modern communication networks. The objectives of congestion control are two-fold: to utilize as much the available capacity of the network as possible without causing severe congestion within the network, and to ensure some form of fairness among the users. Since the seminal work by Kelly [1], it is clear that both of these objectives can be mapped to a global optimization problem that maximizes the total system utility, where different fairness objectives can be achieved by appropriately choosing the utility functions. Congestion control can then be viewed as a distributed iterative solution to the global optimization problem [1, 2, 3, 4].

Significant advances in the understanding of congestion control have been made under this optimization framework for congestion control (see [5] for a good survey). The results can be roughly categorized into two groups. In the first body of work, it is assumed that the number of users in the network is fixed and each user has infinite data to transfer. This research focuses on developing distributed iterative algorithms that converge to the fair rate allocation, which corresponds to the solution of the global optimization problem. Various issues have been addressed in this body of work, including global convergence of the congestion control algorithm, local stability (in the sense of Lyapunov) of the equilibrium rate allocation, the impact of feedback delay and random noise, and the asymptotic behavior of the system when the number of users is large.

The second body of work studies a network with *random dynamic arrivals and departures of the users*. This research studies the *stability region* of the system employing congestion control. Here, by *stability*, we mean that the number of users in the system and the queue lengths at each link in the network remain finite. The *stability region* of the system under a given congestion control algorithm is the set of offered loads under which the system is stable. This body of work typically assumes that, whenever the number of users in the system changes, the data rates of the users are adjusted *instantaneously* to the *optimal (and fair) rate allocation* computed

by the global optimization problem. This model essentially assumes a *time-scale separation*, i.e., the time scale of the arrivals and departures of the users is much slower than that of the dynamics determined by the congestion control algorithms derived in the first body of work. It has been shown that, for a large class of utility functions and fairness objectives, the *largest possible* stability region can be achieved by allocating data rates fairly according to this time-scale separation assumption [6, 7, 8, 9]. This result is important as it tells us that "fairness" is not merely an *aesthetic* property, but it actually has a strong global *performance* implication, i.e., in achieving the *largest possible* stability region.

However, for a large network like today's Internet, with the continual arrivals and departures of the users, the number of users in the system changes constantly. There will rarely be an extended period of time when the number of users in the system is fixed. Hence, the iterative congestion control algorithm in the first body of work may never have the chance to converge to an optimal (and fair) rate allocation. Therefore, the time-scale separation assumption used in the second body of work, albeit technically convenient, rarely holds in practice.

In this paper, we study the stability region of congestion control without requiring such a time-scale separation assumption. We will show that, *even when we remove the time-scale separation assumption*, the *largest possible* stability region can still be achieved by a large class of congestion control algorithms that are derived from the optimization framework. Hence, our result reinforces the performance benefit of congestion control in a stronger sense than previous works.

The rest of the paper is structured as follows. In Section 2, we present the system model and review some relevant results in the literature. Our main result is presented in Section 3, and the proof is given in Section 4. Then we conclude.

## 2   The System Model and Related Results

In this section, we describe our system model and review certain related works. We consider a network with $L$ links and $S$ classes of users. The capacity of each link $l$ is $R^l$. Users of each

class $s$ have one path through the network. Let $H_s^l = 1$, if the path of users of class $s$ uses link $l$, and $H_s^l = 0$, otherwise. Let $x_s$ denote the rate at which each user of class $s$ sends data into the network, and let $U_s(x_s)$ be the utility received by the user of class $s$ when it sends data at rate $x_s$. The utility function $U_s(\cdot)$ characterizes the "satisfaction level" of a user of class $s$ when it sends data at a certain rate, and as we will soon discuss, it also corresponds to a certain fairness objective. As is typically assumed in the literature, we assume that each user of class $s$ has a maximum data-rate limit of $M_s$, and the utility function $U_s(\cdot)$ is increasing, strictly concave, and twice continuously differentiable on $(0, M_s]$ [2].

Let $n_s, s = 1, ..., S$ denote the number of users of class $s$ that are in the system. Let $\vec{n} = [n_1, ..., n_S]$ and $\vec{x} = [x_1, ..., x_S]$. Congestion control can then be formulated as the following global optimization problem [1]:

$$\max_{\vec{x}:0 \leq x_s \leq M_s, s=1,...,S} \sum_{s=1}^{S} n_s U_s(x_s) \tag{1}$$

$$\text{subject to} \quad \sum_{s=1}^{S} H_s^l n_s x_s \leq R^l \quad \text{for all } l = 1, ..., L.$$

## 2.1 Fairness

It has been well known that fairness objectives can be achieved by appropriately choosing the utility functions [7]. For example, utility functions of the form

$$U_s(x_s) = w_s \log x_s \tag{2}$$

correspond to *weighted proportional fairness*, where $w_s, s = 1, ..., S$ are the weights. A more general form of the utility function is

$$U_s(x_s) = w_s \frac{x_s^{1-\beta}}{1 - \beta}, \text{ for some } \beta > 0 \text{ and } \beta \neq 1. \tag{3}$$

Maximizing the total system utility will correspond to *maximizing weighted throughput* as $\beta \to 0$, *weighted proportional fairness* as $\beta \to 1$, and *max-min fainess* as $\beta \to \infty$.

## 2.2 Convergence

We first assume that $\vec{n}$, the number of users in the system, is fixed and each user has an infinite backlog to transfer. We associate an implicit cost $q^l$ with each link $l$ and let $\vec{q} = [q^1, ..., q^L]$. The following iterative algorithm, commonly referred to as the "dual solution" in the congestion control literature, can solve problem (1) with an appropriate choice of the step-size.

**Algorithm $\mathcal{A}$:**

At each time instant $t$,

- The data rate of each user of class $s$ is determined by:

$$x_s(t) = \operatorname*{argmax}_{0 \leq x_s \leq M_s} U_s(x_s) - x_s \sum_{l=1}^{L} H_s^l q^l(t). \tag{4}$$

- The implicit cost at each link $l$ is updated by:

$$q^l(t+1) = \left[ q^l(t) + \alpha_l \left( \sum_{s=1}^{S} H_s^l n_s x_s(t) - R^l \right) \right]^+, \tag{5}$$

where $[\cdot]^+$ denotes the projection to $[0, \infty)$ and $\alpha_l$ is a positive step-size for each link $l$.

The following proposition was shown in [2] with slightly different notation.

**Proposition 1** *Assume that the number of users in the system is fixed. Further, assume that the curvatures of $U_s(\cdot)$ are bounded away from zero on $(0, M_s]$, i.e., there exists a positive number $\gamma_s$ for each class $s$ such that*

$$-U_s''(x_s) \geq \gamma_s > 0 \text{ for all } x_s \in (0, M_s]. \tag{6}$$

*Let $\vec{x}^*$ denote the optimal solution to problem (1). Let $\mathcal{S} = \max_l \sum_{s=1}^{S} H_s^l n_s$ denote the maximum number of **users** using any link, and let $\mathcal{L} = \max_s \sum_{l=1}^{L} H_s^l$ denote the maximum number of links used by any **user**. If*

$$\max_l \alpha_l \leq \frac{2}{\mathcal{S}\mathcal{L}} \min_s \gamma_s, \tag{7}$$

*then Algorithm $\mathcal{A}$ converges, i.e., $\vec{x}(t) \to \vec{x}^*$ as $t \to \infty$.*

## 2.3   Stability Region

We now turn to the case when the number of users in the system changes dynamically. In this case, we will study the *stability region* of system. Here, by *stability*, we mean that the number of users in the system and the queue lengths at each link in the network remain finite. To be precise, we assume that users of class $s$ arrive to the network according to a Poisson process with rate $\lambda_s$ and that each user brings with it a file for transfer whose size is exponentially distributed with mean $1/\mu_s$. The load brought by users of class $s$ is then $\rho_s = \lambda_s/\mu_s$. Let $\vec{\rho} = [\rho_1, ..., \rho_S]$. Let $n_s(t)$ denote the number of users of class $s$ that are in the system at time $t$ and let $\vec{n}(t) = [n_1(t), ..., n_S(t)]$. We assume that the rate allocation for users of the same class is identical. Let $x_s(t)$ denote the rate of users of class $s$ at time $t$ and let $\vec{x}(t) = [x_1(t), ..., x_S(t)]$. In the rate assignment models that follow, the evolution of $\vec{n}(t)$ will be governed by a Markov process. Its transition rates are given by:

$$n_s(t) \rightarrow n_s(t) + 1, \qquad \text{with rate } \lambda_s,$$

$$n_s(t) \rightarrow n_s(t) - 1, \qquad \text{with rate } \mu_s x_s(t) n_s(t) \text{ if } n_s(t) > 0.$$

We say that the above system is *stable* [10] if

$$\limsup_{t \to \infty} \frac{1}{t} \int_0^t \mathbf{1}_{\{\sum\limits_{s=1}^{S} n_s(t) + \sum\limits_{l=1}^{L} q^l(t) > M\}} dt \to 0, \text{ as } M \to \infty.$$

The *stability region* $\Theta$ of the system under a given congestion control algorithm is the set of offered loads $\vec{\rho}$ such that the system is stable for any $\vec{\rho} \in \Theta$.

Past works on the stability region of congestion control typically adopt the following *time-scale separation assumption*:

**The Time-Scale Separation Assumption:**

- The data rates $\vec{x}(t)$ of the users at each time instant $t$ are adjusted *instantaneously* to the optimal rate allocation computed by the global optimization problem (1) with $\vec{n} = \vec{n}(t)$.

We refer to a congestion controller that allocates data rates according to the above time-scale separation assumption as the *perfect congestion controller*.

We say that the stability region achieved by a congestion controller is *the largest possible* when the following holds: for any offered load, if this congestion controller cannot stabilize the system, no other congestion controller can. Note that the capacity constraint determines an *upper bound* on the stability region achieved by *any* congestion controller, i.e.

$$\Theta \subset \Theta_0 \triangleq \left\{ \vec{\rho} \mid \sum_{s=1}^{S} H_s^l \rho_s \leq R^l \text{ for all } l \right\}. \tag{8}$$

The next proposition from [7] shows that the stability region achieved by the *perfect congestion controller* is indeed *the largest possible* found on the right hand side of (8).

**Proposition 2** *Under the time-scale separation assumption, if the utility functions are of the form in (2) or (3) for some $\beta > 0$, then for any offered load $\vec{\rho}$ that resides strictly inside $\Theta_0$, the Markov process $\vec{n}(t)$ is positive recurrent and hence,*

$$\limsup_{t \to \infty} \frac{1}{t} \int_0^t \mathbf{1}_{\{\sum_{s=1}^{S} n_s(t) > M\}} \, dt \to 0, \ \ as \ M \to \infty.$$

# 3   Stability Region of Congestion Control Without the Time-Scale Separation Assumption

As discussed earlier in the Introduction, the time-scale separation assumption rarely holds in reality. In typical networks, users arrive and depart constantly. Hence, the data rates of the users employing a congestion control algorithm such as algorithm $\mathcal{A}$ may never be able to converge. Further, note that the step-size condition (7) in Proposition 1 becomes more stringent as the number of users in the system increases. As the offered load $\vec{\rho}$ approaches the boundary of the stability region $\Theta_0$, the number of users in the system will approach infinity. Hence, given a chosen set of step-sizes, algorithm $\mathcal{A}$ will fail to converge when the offered load is close to the boundary of $\Theta_0$. The time-scale separation assumption will not hold in this case either.

In this section, we will present a new result on the stability region of congestion control *without this time-scale separation assumption.* We first describe some more details of the dynamics of the system. We assume that time is divided into slots of length $T$, and that the implicit costs at the links are updated only at the end of each time slot. However, users may arrive and depart in the middle of a time slot. Let $\vec{q}(kT)$ denote the implicit costs at time slot $k$. Unlike the case in Proposition 2, we now let the rate allocation $\vec{x}(t)$ be determined by the current implicit costs. We assume that the utility function is of the form (2) or (3). Then, by solving (4), the data rate of users of class $s$ is given by

$$x_s(t) = x_s(kT) = \min\left\{\left(\frac{w_s}{\sum_{l=1}^{L} H_s^l q^l(kT)}\right)^{1/\beta}, M_s\right\}, \quad \text{for } kT \le t < (k+1)T \tag{9}$$

(use $\beta = 1$ when the utility functions are of the form (2)). At the end of each time slot, the implicit costs are updated by

$$q^l((k+1)T) = \left[q^l(kT) + \alpha_l\left(\sum_{s=1}^{S} H_s^l \int_{kT}^{(k+1)T} n_s(t)x_s(kT)dt - TR^l\right)\right]^+. \tag{10}$$

The following proposition shows that, even when the time-scale separation assumption is removed, the above congestion control algorithm can still achieve *the largest possible* stability region. The proof is given in Section 4.

**Proposition 3** *Assume that the utility functions are of the form in (2) or (3) for some $\beta > 1$, and that the data rates of the users are controlled by (9). Let $\bar{\mathcal{S}} = \max_l \sum_{s=1}^{S} H_s^l$ denote the maximum number of **classes** using any link, and let $\bar{\mathcal{L}} = \max_s \sum_{l=1}^{L} H_s^l$ denote the maximum number of links used by any **class**, If*

$$\max_l \alpha_l \le \frac{1}{T\bar{\mathcal{S}}\bar{\mathcal{L}}} \frac{2^\beta - 1}{16} \min_s \frac{w_s}{\rho_s M_s^\beta} \tag{11}$$

*(use $\beta = 1$ if the utility functions are of the form in (2)), then for any offered load $\vec{\rho}$ that resides strictly inside $\Theta_0$, the system described by the Markov process $[\vec{n}(kT), \vec{q}(kT)]$ is stable.*

Several remarks are in order: Firstly, no time-scale separation assumption is required in Proposition 3. Hence, we do not require that the data rates of the users converge. Secondly, a step-size rule that is *independent* of the instantaneous number of users in the system is provided in (11) (note the difference between $\mathcal{S}$ and $\bar{\mathcal{S}}$). Given our discussion at the beginning of this section, it is quite remarkable that we do not need to reduce the step-sizes even when the offered load is close to the boundary of the stability region. In fact, since the set $\Theta_0$ is bounded, the step-sizes can be chosen independently of the offered load. The step-size rule (11) is dependent on $M_s$, the maximum data rate of users belonging to class $s$. This dependence is not surprising. Since the utility functions are of the forms in (2) or (3), we have,

$$U_s''(x_s) = -\beta \frac{w_s}{x_s^{\beta+1}}.$$

Hence, the minimum curvature of $U_s(\cdot)$ is

$$\gamma_s = \frac{\beta w_s}{M_s^{\beta+1}}.$$

Let $\tilde{n}_s = \rho_s/M_s$, which can be interpreted as the average number of users of class $s$ in a (fictitious) $M/M/\infty/\infty$ system *where each user of class $s$ is served at its maximum data rate $M_s$*. The step-size condition (11) then becomes

$$\max_l \alpha_l \leq \frac{1}{T\bar{\mathcal{S}}\bar{\mathcal{L}}} \frac{2^\beta - 1}{16\beta} \min_s \frac{\gamma_s}{\tilde{n}_s},$$

which is comparable to (7). However, note that $\tilde{n}_s$ is *quite different* from $\mathbf{E}[n_s(t)]$, the average number of users of class $s$ in the *real* system. Again, as $\tilde{n}_s$ is always bounded, the step-sizes can be chosen independently of the offered load.

# 4  Proof of Proposition 3

Define

$$\mathcal{V}(\vec{n}, \vec{q}) = V_n(\vec{n}) + V_q(\vec{q}),$$

where

$$V_n(\vec{n}) = \frac{1}{(1+\epsilon)^\beta} \sum_{s=1}^{S} \frac{w_s n_s^{\beta+1}}{(1+\beta)\mu_s \rho_s^\beta}, \quad V_q(\vec{q}) = \sum_{l=1}^{L} \frac{(q^l)^2}{2\alpha_l},$$

and $\epsilon$ is a positive constant in $(0, 1]$ to be chosen later. We shall show that $\mathcal{V}(\cdot, \cdot)$ is a Lyapunov function of the system. We begin with a few lemmas. The first lemma bounds the change in $V_n(\cdot)$.

**Lemma 4**

$$\mathbf{E}[V_n(\vec{n}((k+1)T) - V_n(\vec{n}(kT))|\vec{n}(kT), \vec{q}(kT)]$$

$$\leq -\epsilon \sum_{s=1}^{S} E_0(s) \int_{kT}^{(k+1)T} \mathbf{E}[n_s^\beta(t)|\vec{n}(kT), \vec{q}(kT)]dt$$

$$+ \sum_{s=1}^{S} \left[ \sum_{l=1}^{L} H_s^l q^l(kT) \right] \left[ (1+\epsilon)\rho_s T - \int_{t=kT}^{(k+1)T} \mathbf{E}[n_s(t)x_s(t)|\vec{n}(kT), \vec{q}(kT)]dt \right]$$

$$- \sum_{s=1}^{S} \frac{2^\beta - 1}{8(1+\epsilon)} \frac{w_s}{\rho_s M_s^\beta} \int_{kT}^{(k+1)T} \mathbf{E}[n_s^2(t)x_s^2(t)|\vec{n}(kT), \vec{q}(kT)]dt$$

$$+ E_1, \tag{12}$$

where $E_0(s)$ and $E_1$ are finite positive constants.

**Proof:** Over a small time interval $\delta t$, we have

$$\mathbf{E}\left[n_s^{\beta+1}(t+\delta t) - n_s^{\beta+1}(t)|\vec{n}(t), \vec{q}(t)\right]$$

$$= [(n_s(t)+1)^{\beta+1} - n_s^{\beta+1}(t)]\lambda_s \delta t + [(n_s(t)-1)^{\beta+1} - n_s^{\beta+1}(t)]\mu_s n_s(t)x_s(t)\delta t + o(\delta t).$$

By the Mean-Value Theorem,

$$(n + \Delta n)^{\beta+1} - n^{\beta+1} = (\beta+1)n^\beta \Delta n + \frac{\beta(\beta+1)}{2}(n + \nu\Delta n)^{\beta-1}(\Delta n)^2$$

for some $\nu \in (0, 1)$. Hence, letting $\Delta n = \pm 1$, we have

$$\mathbf{E}\left[n_s^{\beta+1}(t+\delta t) - n_s^{\beta+1}(t)|\vec{n}(t), \vec{q}(t)\right]$$

$$\leq (\beta+1)n_s^\beta(t)[\lambda_s \delta t - \mu_s n_s(t)x_s(t)\delta t]$$

$$+ 2^{\beta-2}\beta(\beta+1)n_s^{\beta-1}(t)\left[\lambda_s \delta t + \mu_s n_s(t)x_s(t)\delta t\right] + N_1(s)\delta t + o(\delta t)$$

10

for some positive constant $N_1(s)$. We then have,

$$\frac{\mathbf{E}[V_n(\vec{n}(t+\delta t)) - V_n(\vec{n}(t))|\vec{n}(t), \vec{q}(t)]}{\delta t}$$

$$\leq \frac{1}{(1+\epsilon)^\beta} \sum_{s=1}^{S} \left\{ \frac{w_s n_s^\beta(t)}{\mu_s \rho_s^\beta} [\lambda_s - \mu_s n_s(t) x_s(t)] \right.$$

$$\left. + \frac{\beta 2^{\beta-2} w_s n_s^{\beta-1}(t)}{\mu_s \rho_s^\beta} [\lambda_s + \mu_s n_s(t) x_s(t)] + N_1(s) \right\} + o(1)$$

$$= \frac{1}{(1+\epsilon)^\beta} \sum_{s=1}^{S} \left\{ \frac{w_s n_s^\beta(t)}{\rho_s^\beta} [\rho_s - n_s(t) x_s(t)] \right.$$

$$\left. + \frac{\beta 2^{\beta-2} w_s n_s^{\beta-1}(t)}{\rho_s^\beta} [\rho_s + n_s(t) x_s(t)] + N_1(s) \right\} + o(1)$$

$$= -\frac{\epsilon}{(1+\epsilon)^\beta} \sum_{s=1}^{S} \frac{w_s n_s^\beta(t)}{\rho_s^{\beta-1}} + \frac{1}{(1+\epsilon)^\beta} \sum_{s=1}^{S} \left\{ \frac{w_s n_s^\beta(t)}{\rho_s^\beta} [(1+\epsilon)\rho_s - n_s(t) x_s(t)] \right.$$

$$\left. + \frac{\beta 2^{\beta-2} w_s n_s^{\beta-1}(t)}{\rho_s^\beta} [\rho_s + n_s(t) x_s(t)] + N_1(s) \right\} + o(1) \qquad (13)$$

$$\leq -\epsilon \sum_{s=1}^{S} N_0(s) n_s^\beta(t) + \sum_{s=1}^{S} \left\{ \left[ \sum_{l=1}^{L} H_s^l q^l(t) \right] [(1+\epsilon)\rho_s - n_s(t) x_s(t)] \right.$$

$$+ \left[ \frac{w_s}{x_s^\beta(t)} - \sum_{l=1}^{L} H_s^l q^l(t) \right] [(1+\epsilon)\rho_s - n_s(t) x_s(t)] \qquad (14)$$

$$+ w_s \left[ \frac{n_s^\beta(t)}{((1+\epsilon)\rho_s)^\beta} - \frac{1}{x_s^\beta(t)} \right] [(1+\epsilon)\rho_s - n_s(t) x_s(t)] \qquad (15)$$

$$+ \frac{\beta 2^{\beta-2}}{(1+\epsilon)^\beta} \frac{w_s n_s^{\beta-1}(t)}{\rho_s^\beta} [\rho_s + n_s(t) x_s(t)] \qquad (16)$$

$$\left. + \frac{N_1(s)}{(1+\epsilon)^\beta} \right\} + o(1),$$

where

$$N_0(s) = \frac{1}{(1+\epsilon)^\beta} \frac{w_s}{\rho_s^{\beta-1}}.$$

We shall bound the three terms (14-16). By (9),

$$\frac{w_s}{x_s^\beta(t)} = \max \left\{ \sum_{l=1}^{L} H_s^l q^l(t), \frac{w_s}{M_s^\beta} \right\}.$$

Hence, the term (14) can be bounded by

$$\left[ \frac{w_s}{x_s^\beta(t)} - \sum_{l=1}^{L} H_s^l q^l(t) \right] [(1+\epsilon)\rho_s - n_s(t) x_s(t)]$$

11

$$\leq \left[\frac{w_s}{x_s^\beta(t)} - \sum_{l=1}^{L} H_s^l q^l(t)\right](1+\epsilon)\rho_s$$

$$\leq \left[\frac{w_s}{M_s^\beta} - \sum_{l=1}^{L} H_s^l q^l(t)\right]^+ (1+\epsilon)\rho_s$$

$$\leq N_2(s) \triangleq \frac{(1+\epsilon)w_s\rho_s}{M_s^\beta}. \tag{17}$$

Let $(A)$ and $(B)$ denote the terms (15) and (16), respectively. Note that

$$(A) = w_s \left[\frac{n_s^\beta(t)}{((1+\epsilon)\rho_s)^\beta} - \frac{1}{x_s^\beta(t)}\right][(1+\epsilon)\rho_s - n_s(t)x_s(t)]$$

$$= -w_s \frac{[(1+\epsilon)\rho_s - n_s(t)x_s(t)]\left[((1+\epsilon)\rho_s)^\beta - n_s^\beta(t)x_s^\beta(t)\right]}{((1+\epsilon)\rho_s)^\beta x_s^\beta(t)} \leq 0. \tag{18}$$

If $n_s(t)x_s(t) \geq 2(1+\epsilon)\rho_s$, then

$$[(1+\epsilon)\rho_s - n_s(t)x_s(t)]\left[((1+\epsilon)\rho_s)^\beta - n_s^\beta(t)x_s^\beta(t)\right]$$

$$\geq \left[\frac{n_s(t)x_s(t)}{2}\right]\left[\frac{2^\beta - 1}{2^\beta}n_s^\beta(t)x_s^\beta(t)\right]. \tag{19}$$

Hence,

$$(A) \leq -\frac{2^\beta - 1}{2^{\beta+1}}\frac{w_s n_s^{\beta+1}x_s(t)}{((1+\epsilon)\rho_s)^\beta},$$

and

$$(B) \leq \frac{\beta 2^{\beta-1}}{(1+\epsilon)^\beta}\frac{w_s n_s^\beta(t)x_s(t)}{\rho_s^\beta}.$$

Since

$$\beta 2^{\beta-1}n_s^\beta(t) \leq \frac{2^\beta - 1}{2^{\beta+2}}n_s^{\beta+1}(t) + N_3(s)$$

for some positive constant $N_3(s)$, we have,

$$(B) \leq -\frac{(A)}{2} + \frac{w_s N_3(s)x_s(t)}{((1+\epsilon)\rho_s)^\beta}$$

$$\leq -\frac{(A)}{2} + N_4(s),$$

where

$$N_4(s) = \frac{w_s N_3(s)M_s}{((1+\epsilon)\rho_s)^\beta}.$$

12

On the other hand, if $n_s(t)x_s(t) < 2(1 + \epsilon)\rho_s \le 4\rho_s$, then

$$
\begin{aligned}
(B) &\le \frac{5\beta 2^{\beta-2}}{(1+\epsilon)^\beta} \frac{w_s n_s^{\beta-1}(t)}{\rho_s^{\beta-1}} \\
&= N_5(s) n_s^{\beta-1}(t),
\end{aligned}
$$

where

$$
N_5(s) = \frac{5\beta 2^{\beta-2}}{(1+\epsilon)^\beta} \frac{w_s}{\rho_s^{\beta-1}}.
$$

In both cases,

$$
(B) \le -\frac{(A)}{2} + N_5(s) n_s^{\beta-1}(t) + N_4(s). \tag{20}
$$

Substituting (17) and (20) back to (14-16), we have,

$$
\begin{aligned}
&\frac{\mathbf{E}[V_n(\vec{n}(t + \delta t)) - V_n(\vec{n}(t)) | \vec{n}(t), \vec{q}(t)]}{\delta t} \\
&\le -\sum_{s=1}^{S} \left[ \epsilon N_0(s) n_s^\beta(t) - N_5(s) n_s^{\beta-1}(t) \right] \\
&\quad + \sum_{s=1}^{S} \left[ \sum_{l=1}^{L} H_s^l q^l(t) \right] [(1+\epsilon)\rho_s - n_s(t)x_s(t)] \\
&\quad + \sum_{s=1}^{S} \frac{(A)}{2} + \sum_{s=1}^{S} [N_1(s) + N_2(s) + N_4(s)] + o(1). \tag{21}
\end{aligned}
$$

We shall use (18) and (19) again to bound $(A)/2$. Since $x_s(t) \le M_s$, if $n_s(t)x_s(t) \ge 2(1 + \epsilon)\rho_s$, we have,

$$
\begin{aligned}
\frac{(A)}{2} &\le -w_s \frac{2^\beta - 1}{2^{\beta+2}} \frac{n_s^{\beta+1}(t)x_s^{\beta+1}(t)}{((1+\epsilon)\rho_s)^\beta M_s^\beta} \\
&\le -w_s \frac{2^\beta - 1}{2^{\beta+2}} \frac{2^{\beta-1}}{(1+\epsilon)\rho_s} \frac{n_s^2(t)x_s^2(t)}{M_s^\beta} \\
&\le -w_s \frac{2^\beta - 1}{8(1+\epsilon)} \frac{n_s^2(t)x_s^2(t)}{\rho_s M_s^\beta}.
\end{aligned} \tag{22}
$$

On the other hand, if $n_s(t)x_s(t) < 2(1 + \epsilon)\rho_s$, we still have $(A)/2 \le 0$. Hence, in both cases,

$$
\frac{(A)}{2} \le -w_s \frac{2^\beta - 1}{8(1+\epsilon)} \frac{n_s^2(t)x_s^2(t)}{\rho_s M_s^\beta} + N_6(s), \tag{23}
$$

13

where

$$N_6(s) = w_s \frac{2^\beta - 1}{8(1 + \epsilon)} \frac{(2(1 + \epsilon)\rho_s)^2}{\rho_s M_s^\beta}.$$

Further, note that

$$N_5(s)n_s^{\beta-1}(t) \leq \frac{\epsilon N_0(s)}{2} n_s^\beta(t) + N_7(s) \tag{24}$$

for some positive constant $N_7(s)$. Substituting (23) and (24) into (21), and integrating over $[kT, (k+1)T]$, the result (12) then follows with $E_0(s) = N_0(s)/2$ and

$$E_1 = T \sum_{s=1}^{S} [N_1(s) + N_2(s) + N_4(s) + N_6(s) + N_7(s)].$$

*Q.E.D.*

The next lemma bounds the change in $V_q(\cdot)$. For simplicity, we use the following matrix notation. Let $A$ denote the $L \times L$ diagonal matrix whose $l$-th diagonal element is $\alpha_l$. Let $H$ denote the $L \times S$ matrix whose $(l, s)$-element is $H_s^l$. Let $\vec{R} = [R^1, ..., R^l]^{\text{tr}}$, where $[\cdot]^{\text{tr}}$ denotes the transpose. Further let $X_s(t) = n_s(t)x_s(t)$ and let $\vec{X}(t) = [X_1(t), ..., X_S(t)]^{\text{tr}}$. Then

$$V_q(\vec{q}) = \frac{\vec{q}^{\text{tr}} A^{-1} \vec{q}}{2},$$

and the update on the implicit costs (10) can be written as

$$\vec{q}((k+1)T) = \left[ \vec{q}(kT) + A \left( H \int_{kT}^{(k+1)T} \vec{X}(t)dt - \vec{R}T \right) \right]^+. \tag{25}$$

**Lemma 5**

$$\mathbf{E}[V_q(\vec{q}((k+1)T) - V_q(\vec{q}(kT))|\vec{n}(kT), \vec{q}(kT)]$$
$$\leq \vec{q}^{\text{tr}}(kT) \left[ H \int_{kT}^{(k+1)T} \mathbf{E}[\vec{X}(t)|\vec{n}(kT), \vec{q}(kT)]dt - \vec{R}T \right]$$
$$+ T\alpha_{\max} \bar{S} \bar{\mathcal{L}} \sum_{s=1}^{S} \left[ \int_{kT}^{(k+1)T} \mathbf{E}[n_s^2(t)x_s^2(t)|\vec{n}(kT), \vec{q}(kT)]dt \right] + E_2, \tag{26}$$

where $\alpha_{\max} = \max_l \alpha^l$, $\bar{\mathcal{L}}$ and $\bar{S}$ are defined as in Proposition 3, and $E_2$ is a finite positive constant.

**Proof:** By (25),

$$
\begin{aligned}
V_q(&\vec{q}((k+1)T)) - V_q(\vec{q}(kT)) \\
\leq\ & \vec{q}^{\,\mathrm{tr}}(kT)\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt - \vec{R}T\right] \\
&+\frac{1}{2}\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt - \vec{R}T\right]^{\mathrm{tr}} A\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt - \vec{R}T\right] \\
\leq\ & \vec{q}^{\,\mathrm{tr}}(kT)\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt - \vec{R}T\right] \\
&+\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt\right]^{\mathrm{tr}} A\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt\right] + T^2\vec{R}^{\,\mathrm{tr}}A\vec{R}.
\end{aligned}
$$

For the second term, we have,

$$
\begin{aligned}
&\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt\right]^{\mathrm{tr}} A\left[H\int_{kT}^{(k+1)T}\vec{X}(t)dt\right] \\
=\ & \sum_{l=1}^{L}\alpha_l\left[\sum_{s=1}^{S}H_s^l\int_{kT}^{(k+1)T}n_s(t)x_s(t)dt\right]^2 \\
\leq\ & \sum_{l=1}^{L}\alpha_l\left[\sum_{s=1}^{S}H_s^l\right]\left[\sum_{s=1}^{S}H_s^l\left(\int_{kT}^{(k+1)T}n_s(t)x_s(t)dt\right)^2\right] \\
\leq\ & \bar{\mathcal{S}}\sum_{l=1}^{L}\alpha_l\left[\sum_{s=1}^{S}H_s^l\left(\int_{kT}^{(k+1)T}n_s(t)x_s(t)dt\right)^2\right] \\
\leq\ & T\bar{\mathcal{S}}\sum_{l=1}^{L}\alpha_l\sum_{s=1}^{S}H_s^l\int_{kT}^{(k+1)T}n_s^2(t)x_s^2(t)dt \\
=\ & T\bar{\mathcal{S}}\sum_{s=1}^{S}\left[\int_{kT}^{(k+1)T}n_s^2(t)x_s^2(t)dt\right]\left[\sum_{l=1}^{L}\alpha_lH_s^l\right] \\
\leq\ & T\alpha_{\max}\bar{\mathcal{S}}\bar{\mathcal{L}}\sum_{s=1}^{S}\left[\int_{kT}^{(k+1)T}n_s^2(t)x_s^2(t)dt\right]. 
\end{aligned}
\tag{27}
$$

Letting $E_2 = T^2\vec{R}^{\,\mathrm{tr}}A\vec{R}$, the result (26) then follows. $\hspace{3cm}$ Q.E.D.

**Proof of Proposition 3 :** Adding (12) to (26), and noting that

$$
\sum_{s=1}^{S}\left\{\left[\sum_{l=1}^{L}H_s^l q^l(kT)\right]\int_{kT}^{(k+1)T}\mathbf{E}[n_s(t)x_s(t)|\vec{n}(kT),\vec{q}(kT)]dt\right\}
$$

15

$$= \sum_{l=1}^{L} q^l(kT) \sum_{s=1}^{S} H_s^l \int_{kT}^{(k+1)T} \mathbf{E}[n_s(t)x_s(t)|\vec{n}(kT), \vec{q}(kT)]dt$$

$$= \vec{q}^{\,\text{tr}}(kT)\left[H\int_{kT}^{(k+1)T} \mathbf{E}[\vec{X}(t)|\vec{n}(kT), \vec{q}(kT)]dt\right],$$

we have,

$$\mathbf{E}[\mathcal{V}(\vec{n}((k+1)T), \vec{q}((k+1)T)) - \mathcal{V}(\vec{n}(kT), \vec{q}(kT))|\vec{n}(kT), \vec{q}(kT)]$$

$$\leq -\epsilon \sum_{s=1}^{S} E_0(s) \int_{kT}^{(k+1)T} \mathbf{E}[n_s^\beta(t)|\vec{n}(kT), \vec{q}(kT)]dt$$

$$+ \sum_{s=1}^{S}\left[\sum_{l=1}^{L} H_s^l q^l(kT)\right](1+\epsilon)\rho_s T - \vec{q}^{\,\text{tr}}(kT)\vec{R}T$$

$$- \sum_{s=1}^{S}\left[\frac{2^\beta - 1}{8(1+\epsilon)}\frac{w_s}{\rho_s M_s^\beta} - T\alpha_{\max}\bar{\mathcal{S}}\bar{\mathcal{L}}\right]$$

$$\times \int_{kT}^{(k+1)T} \mathbf{E}[n_s^2(t)x_s^2(t)|\vec{n}(kT), \vec{q}(kT)]dt \qquad (28)$$

$$+ E_3,$$

where $E_3 = E_1 + E_2$. If (11) is satisfied, then the product term in (28) is negative. Hence, by some rearrangement of the order of the summations, we have,

$$\mathbf{E}[\mathcal{V}(\vec{n}((k+1)T), \vec{q}((k+1)T)) - \mathcal{V}(\vec{n}(kT), \vec{q}(kT))|\vec{n}(kT), \vec{q}(kT)]$$

$$\leq -\epsilon \sum_{s=1}^{S} E_0(s) \int_{kT}^{(k+1)T} \mathbf{E}[n_s^\beta(t)|\vec{n}(kT), \vec{q}(kT)]dt + T\vec{q}^{\,\text{tr}}(kT)\left[(1+\epsilon)H\vec{\rho} - \vec{R}\right] + E_3.$$

By assumption, $\vec{\rho}$ lies strictly inside $\Theta_0$. Hence, there exists some $\epsilon \in (0, 1]$ such that $(1+2\epsilon)H\vec{\rho} \leq \vec{R}$. Use this value of $\epsilon$ in the definition of $\mathcal{V}(\cdot, \cdot)$. we then have,

$$\mathbf{E}[\mathcal{V}(\vec{n}((k+1)T), \vec{q}((k+1)T)) - \mathcal{V}(\vec{n}(kT), \vec{q}(kT))|\vec{n}(kT), \vec{q}(kT)]$$

$$\leq -\epsilon \sum_{s=1}^{S} E_0(s) \int_{kT}^{(k+1)T} \mathbf{E}[n_s^\beta(t)|\vec{n}(kT), \vec{q}(kT)]dt - \epsilon T\vec{q}^{\,\text{tr}}(kT)H\vec{\rho} + E_3$$

$$\leq -\epsilon'\left[\sum_{s=1}^{S} n_s^\beta(kT) + \sum_{l=1}^{L} q^l(kT)\right] + E_3$$

for some $\epsilon' > 0$. By Theorem 2 of [10], the result then follows. $\hspace{2cm}$ Q.E.D.

*Remark:* This proof will not work for $\beta < 1$, in which case the relationship (22) will fail to hold. (We need (22) to cancel the second term in (26) of the change in $V_q(\cdot)$.) We have not been able to either prove or disprove our result for $\beta < 1$. We could have resorted to the fluid limit technique of [11]. However, the difficulty in applying the technique of [11] is that the fluid limit of our system is not well defined whenever $\sum_{l=1}^{L} H_s^l q^l(t) = 0$ for some class $s$, which also corresponds to the case when the second term in (26) is large. We will leave the case $\beta < 1$ for future work.

# 5    Conclusion

In this paper, we have studied the stability region of a network employing congestion control algorithms derived from an optimization framework. We have removed the *time-scale separation assumption* typical in other related works, and established that the *largest possible* stability region can be achieved by a large class of congestion control algorithms (i.e., the so-called "dual solutions") derived from the optimization framework. Our result provides new insights on the performance implication of congestion control, and on the choices of the parameters of the congestion controller.

Several directions for future work are possible. Firstly, it would be interesting to see whether our main result holds for the so-called "primal solutions" in the literature [1]. Secondly, we have assumed a Markovian model in this paper. We plan to extend our result to more general user arrival and departure processes. Our result can also be extended to other forms of utility functions [9]. Thirdly, we plan to study the impact of feedback delay. We expect that our main result (Proposition 3) would hold even in the presence of feedback delay, provided that the step-sizes are appropriately chosen. Finally, the extension to the case with multipath routing would also be interesting.

# References

[1] F. P. Kelly, A. Maulloo, and D. Tan, "Rate Control in Communication Networks: Shadow Prices, Proportional Fairness and Stability," *Journal of the Operational Research Society*, vol. 49, pp. 237–252, 1998.

[2] S. H. Low and D. E. Lapsley, "Optimization Flow Control–I: Basic Algorithm and Convergence," *IEEE/ACM Transactions on Networking*, vol. 7, no. 6, pp. 861–874, December 1999.

[3] S. Kunniyur and R. Srikant, "End-to-End Congestion Control Schemes: Utility Functions, Random Losses and ECN Marks," *IEEE/ACM Transactions on Networking*, vol. 11, no. 5, pp. 689–702, October 2003.

[4] H. Yaiche, R. Mazumdar, and C. Rosenberg, "A Game Theoretic Framework for Bandwidth Allocation and Pricing in Broadband Networks," *IEEE/ACM Transactions on Networking*, vol. 8, no. 5, pp. 667–678, Oct. 2000.

[5] S. H. Low and R. Srikant, "A Mathematical Framework for Designing a Low-Loss Low-Delay Internet," *Network and Spatial Economics*, vol. 4, no. 1, pp. 75–102, March 2004.

[6] G. De Veciana, T. J. Lee, and T. Konstantopoulos, "Stability and Performance Analysis of Networks Supporting Elastic Services," *IEEE/ACM Transactions on Networking*, vol. 9, no. 1, pp. 2–14, February 2001.

[7] T. Bonald and L. Massoulie, "Impact of Fairness on Internet Performance," in *Proceedings of ACM Sigmetrics*, Cambridge, MA, June 2001, pp. 82–91.

[8] G. Fayolle, A. L. Fortelle, J. M. Lasgouttes, L. Massoulie, and J. Roberts, "Best Effort Networks: Modeling and Performance Analysis via Large Network Asymptotics," in *Proceedings of IEEE INFOCOM*, Anchorage, Alaska, April 2001.

[9] H. Q. Ye, "Stability of Data Networks Under an Optimization-Based Bandwidth Allocation," *IEEE Transactions on Automatic Control*, vol. 48, no. 7, pp. 1238–1242, July 2003.

[10] M. J. Neely, E. Modiano, and C. E. Rohrs, "Power Allocation and Routing in Multibeam Satellites with Time-Varying Channels," *IEEE/ACM Transactions on Networking*, vol. 11, no. 1, pp. 138–152, February 2003.

[11] J. G. Dai, "On Positive Harris Recurrence of Multiclass Queueing Networks: A Unified Approach via Fluid Limit Models," *Annals of Applied Probability*, vol. 5, pp. 49–77, 1995.