

TCP-Illinois: A Loss and Delay-Based Congestion Control Algorithm for High-Speed Networks

Shao Liu, Tamer Başar and R. Srikant

Abstract— We introduce a new congestion control algorithm for high speed networks, called TCP-Illinois. TCP-Illinois uses packet loss information to determine whether the window size should be increased or decreased, and uses queueing delay information to determine the amount of increment or decrement. TCP-Illinois achieves high throughput, allocates the network resource fairly, and is incentive compatible with standard TCP. We also build a new stochastic matrix model, capturing standard TCP and TCP-Illinois as special cases, and use this model to analyze their fairness properties for both synchronized and unsynchronized backoff behaviors. We finally perform simulations to demonstrate the performance of TCP-Illinois.

Keywords: Congestion Control, TCP, fairness, stability, synchronization

I. INTRODUCTION

TCP-Reno [12], TCP-NewReno [9], and SACK TCP [21] are the standard versions of TCP congestion control protocols currently deployed in the Internet, and they have achieved great success in performing congestion avoidance and control. The key feature of standard TCP is its congestion avoidance phase, which uses the additive increment multiplicative decrement (AIMD) algorithm [11]. Being a window-based algorithm, TCP controls the send rate by maintaining a window size variable W , which limits the number of unacknowledged packets in the network from a single user. This window size is adjusted by the AIMD algorithm in the following manner: W is increased by α/W ($\alpha = 1$ for standard setting) for each ACK, and thus is increased by a constant α/b per round trip time (RTT) if all the packets are acknowledged within an RTT, where b is the number of packets acknowledged by each ACK ($b = 1$ for original TCP, and $b = 2$ for delayed ACK [26]). On the other hand, W is decreased by a fixed proportion βW ($\beta = 1/2$ for standard setting) once some packets are detected to be lost in the last RTT¹. Under this algorithm, senders gently probe the network for spare bandwidth by cautiously increasing their send rates, and sharply reduce their send rates when congestion is detected. Along with other features like slow start, fast recovery, and fast retransmission, TCP

achieves congestion control successfully in the current low speed networks.

However, the current TCP can perform poorly in networks with high bandwidth-delay product (BDP) paths, since the AIMD algorithm, being very conservative, is not designed for large window size flows. First, it takes too long time for a large window size user to recover after a backoff and the bandwidth is not effectively utilized [8]. Second, TCP's time average window size \bar{W} is related with the loss event probability² p in the following manner [23]

$$\bar{W} \approx \sqrt{3/2bp} \quad \text{or} \quad p \approx \frac{3}{2b(\bar{W})^2}. \quad (1)$$

Since TCP interprets all packet losses as congestion signals, \bar{W} is upper bounded by $\sqrt{3/2bp_t}$, where p_t is the transmission error rate [8]. p_t is around 10^{-7} in optical fiber networks, and much higher in other lossy networks, like wireless networks. So TCP, and its AIMD algorithm in particular, should be modified in high bandwidth delay product networks.

Several alternatives to current versions of TCP have been proposed for implementation in high-speed networks. Some require the modification to router algorithms also, like XCP [14], and some modify the sender side only, like HS-TCP [8], Scalable TCP [15], TCP-Westwood [31], H-TCP [17], BIC-TCP [30], TCP Vegas [7], FAST TCP [13] and Compound-TCP [27]. Although each of these has shown advantages over standard TCP in some aspects, none of them have yet provided convincing evidence that they are overwhelmingly better than standard TCP and are suitable for general deployment.

In this paper, we first list some desirable design specifications that a high speed TCP variant should meet, and then introduce the TCP-Illinois algorithm, which uses packet loss information as the *primary* congestion signal to determine the *direction* of window size change (whether window size should be increased or decreased), and uses queueing delay information as the *secondary* congestion signal to adjust the *pace* of window size change (the amount of window size increment or decrement). We then show that TCP-Illinois satisfies all the design requirements we listed, and outperforms standard TCP and many other TCP variants.

To study the fairness, stability, and responsiveness properties of TCP-Illinois, we extend the stochastic matrix model [2]–[4], [16], [25], [29] by allowing window size backoff

All the three authors are with the Department of Electrical and Computer Engineering and Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, 1308 West Main Street, Urbana, IL 61801-2307, USA. Emails: (shaoliu,basar1,rsrikant)@uiuc.edu

Research supported by the NSF ITR Grant CCR 00-85917.

This is an extended version of [19], a paper with the same title, which was presented at First International Conference on Performance Evaluation Methodologies and Tools (VALUETOOLS), October, 2006.

¹Within one RTT, W may decrease multiple times in Reno and can decrease only once in NewReno and SACK.

²All the packet losses within one RTT are regarded as one loss event. Loss event probability is the number of loss events divided by the number of packets sent.

probabilities to be functions of flow arrival rates at congestion events. Our contribution to this modeling technique includes the following: (i) we show that a large class of general AIMD algorithms, including standard TCP and TCP-Illinois, have similar fairness properties, and the fairness properties only depend on the backoff behaviors for these algorithms; (ii) the backoff behavior can be characterized by a function $f(\cdot)$, where $f(\cdot)$ is a user's backoff probability as a function of its send rate: if $f(\cdot) \equiv 1$, the backoffs of all the users are completely synchronized; if $f(\cdot)$ is linear, the backoffs are completely unsynchronized; and in general, the partially-unsynchronized backoff lies in the middle; (iii) $f(\cdot)$ is determined by the number of packets dropped in each congestion event: heavy congestion causes synchronization and light congestion leads to unsynchronized backoff; (iv) the heaviness of congestion (the number of packets dropped in one congestion event as compared to the number of flows) depends on the window size increment of these flows just before the congestion event, and thus a smaller (respectively, larger) increment before a congestion event causes the backoff to be more unsynchronized (respectively, synchronized).

The paper is structured as follows. In section II, we list the requirements for a new version of TCP, compare the existing protocols and point out their shortcomings, and describe our design objective. We next introduce the TCP-Illinois protocol in section III, and study its fairness and stability properties using a new stochastic matrix model for a class of general AIMD algorithms in section IV. We further explore some other properties of TCP-Illinois in section V and provide *ns-2* simulation results in section VI for a comparative study of TCP, HS-TCP and TCP-Illinois.

II. BACKGROUND AND MOTIVATION

As we have mentioned above, several new protocols have been introduced to replace standard TCP in high speed networks. To compare these protocols and to provide insight into the development of an ideal protocol, we list below some requirements that a new protocol should satisfy. This list broadens the list of requirements in [17].

Intra-protocol requirements: The requirements that the protocol should satisfy in a network consisting of a single protocol are the following:

- **Efficiency.** The average throughput for the new protocol should be larger than that of standard TCP in high speed networks.
- **Intra-Protocol Fairness.** Network resources should be fairly allocated to all flows. Fairness here does not necessarily mean that all flows sharing the same link achieve the same throughput. Instead, this means that the new protocol should not be significantly more unfair than the current TCP. For example, under the current TCP, flows with different RTTs achieve similar average window sizes, and their average throughput are inversely proportional to their RTTs. The new protocol should not be significantly more biased against long RTT flows.

- **Responsiveness.** The congestion control algorithm should reach the fair operating point quickly, starting from any initial condition.
- **Heavy Congestion³ Avoidance.** A simple idea to achieve a larger average window size for a given loss event probability is to choose a large value for α (window increment parameter) and a small value for β (window reduction parameter). However, rapid increase and small decrease in window size may cause large number of packets to be dropped during a congestion event, which we call heavy congestion, and thus may lead to some undesirable consequences. First, heavy congestion causes more timeouts and makes TCP enter the slow start phase more often, and causes under-utilization. For example, HS-TCP faces timeouts regularly if SACK is not used. Second, heavy congestion causes synchronization more often, which makes the resource allocation very unfair for large RTT users, as will be discussed later.
- **Router Independence.** The new protocol should work well regardless of router characteristics, like the buffer size at the router, and the queue management algorithm of the router (Droptail or some Active Queue Management (AQM) schemes). With a more advanced router, like with a larger buffer or an AQM support, the new protocol might achieve better performance, but the performance with Droptail and small buffer should also be good.
- **Robustness.** The new protocol should be robust against the noise in congestion signal measurements, especially if this new protocol uses queueing delay as the congestion signal, since queueing delay measurements are typically noisy.

Inter-protocol requirements: The requirements on the protocol when it co-exists in a network with standard TCP are the following:

- **Compatibility.** In low speed networks, the new protocol should achieve a similar rate as that of standard TCP; and in high-speed networks, standard TCP should not suffer significant throughput loss when it coexists with the new protocol.
- **Incentive to switch.** By switching to the new protocol from standard TCP, the users should achieve a higher average throughput in a network that accommodates both protocols.

We now briefly discuss existing TCP variants to see whether they satisfy all these requirements. First, it is impractical to modify routers if the benefit is marginal or can be achieved by sender side modifications, and thus algorithms which need router side modifications, like XCP, are not ideal. Without modifying the router, a sender has only two congestion signals: packet loss and queueing delay. We can

³In our context, heavy congestion means that many packets are dropped when congestion happens. It only concerns the time when congestion happens and it does not necessarily mean that the packet loss probability or the loss event probability is high. In some other papers, it is called heavy synchronization.

thus classify the prior sender-side protocols into one of two classes. Loss-based congestion control algorithms, like HS-TCP and Scalable TCP, use packet loss as primary congestion signal, increase window size for each ACK and decrease window size for packet loss. Loss-based algorithms can be regarded as generalizations of TCP's AIMD algorithm, and we call them general AIMD algorithms, since the only difference from AIMD is that they set different α and β values and allow them to be variables. On the other hand, delay-based congestion control algorithms, like TCP-Vegas and FAST TCP, are fundamentally different from AIMD, as they use queueing delay as the primary congestion signal, increase window size if delay is small and decrease window size if delay is large.

The advantage of delay-based algorithms is that they achieve better average throughput, since they can keep the system around full utilization. As a comparison, the loss-based algorithms purposely generate packet losses and oscillate between full utilization and under utilization. However, existing delay-based algorithms suffer from some inherent weaknesses. First, they are not compatible with standard TCP. TCP-Vegas gets a very small share of the link capacity if competing with TCP-Reno [22], [1]; and FAST TCP yields non-unique equilibrium point if competing with TCP-Reno: the allocation of the bandwidth between FAST and Reno users depend on which users enter the network first [28]. Second, they require the buffer size at the router to be larger than a specified value and this value increases with the number of users N . Both Vegas and FAST control the number of packets queued in the router for each flow, and this number cannot be too small. The requirement for the router buffer is thus N times this number. For a fixed buffer size, there is an upper bound on N for Vegas or FAST to work efficiently. Finally, the performance of these delay-based algorithms deteriorates if the delay measurements are noisy [6], [20], [24].

On the other hand, none of the existing loss-based algorithms satisfy all the requirements either. Scalable TCP sets α proportional to W , but it has been demonstrated to be unfair (see [17], Fig. 2). HS-TCP sets α to be a step-wise increasing function of W , and β a step-wise decreasing function of W , but its convergence speed is very slow (see [17], Fig. 1). H-TCP aims at a faster convergence and better utilization by setting α to be an increasing function of the time elapsed since last backoff and setting β to be such that the link is always around full utilization, even after the backoff. For all the above algorithms, the increase is initially slow, when the window size is small and the network is far from congestion, but becomes fast later, when the window size is large and the network is close to congestion. As a result, the window size curve between two consecutive loss events is convex. This convex nature is not desirable. First, the slow increment in window size when the network is far from congestion is inefficient. For a given β , the convex window curve gets an even smaller average throughput than traditional linear increase, and thus these algorithms have to choose a smaller $\beta < 1/2$, which is not friendly to standard

TCP. Second, the fast increment in window size when the network is close to congestion causes heavy congestion more easily. As we have mentioned before and will further discuss later, heavy congestion causes more frequent timeouts, more synchronized window backoffs, and is more unfair to large RTT users. In summary, the main problem with existing general AIMD algorithms is the convexity of the W curve. An ideal window curve should be concave, which is more efficient and avoids heavy congestion. An objective of our work is to design a general AIMD algorithm which results in a concave window curve.

III. THE TCP-ILLINOIS PROTOCOL

To achieve the concave window curve, we should set α large when far from congestion and set it small when close to congestion. To achieve a better throughput in networks with packet losses not due to congestion and to be fair with standard TCP, we should also set β small when far from congestion and set it large when close to congestion. The difficulty is in judging whether the congestion is imminent or not, since it requires an estimation of the current congestion level. Before congestion (packet loss) really happens, the only congestion indicating information is queueing delay. So our key idea is the following: when the average queueing delay d_a is small, the sender assumes that the congestion is not imminent and sets a large α and small β ; when d_a is large, the sender assumes that the congestion is imminent and sets a small α and large β . As a result, $\alpha = f_1(d_a)$ and $\beta = f_2(d_a)$, where $f_1(\cdot)$ is decreasing and $f_2(\cdot)$ is increasing. Any combination of increasing $f_1(\cdot)$ and decreasing $f_2(\cdot)$ functions results in a concave window curve and therefore, we call such algorithms Concave-AIMD or C-AIMD algorithms.

Note that C-AIMD algorithms use loss to determine the *direction* and use delay to adjust the *pace* of window size change. So loss is the primary congestion signal and delay is the secondary congestion signal. This makes C-AIMD fundamentally different from another recently proposed algorithm, called Compound-TCP [27], which uses both loss and delay information as primary signals (determining *direction* of window size change). As we have mentioned, one problem in using delay to control congestion is that delay cannot be measured accurately since usually the RTT measurements are noisy. If delay determines the *direction* of window size change, noisy RTT measurements could degrade the performance significantly. Our C-AIMD algorithms which use delay only as a *secondary* signal, are much more robust to noise in RTT measurements, as discussed in subsection VI-D.

There are numerous choices for $f_1(\cdot)$ and $f_2(\cdot)$. TCP-Illinois is a special case of C-AIMD algorithms which uses the following choices for $f_1(\cdot)$ and $f_2(\cdot)$:

$$\alpha = f_1(d_a) = \begin{cases} \alpha_{max} & \text{if } d_a \leq d_1 \\ \frac{\kappa_1}{\kappa_2 + d_a} & \text{otherwise.} \end{cases} \quad (2)$$

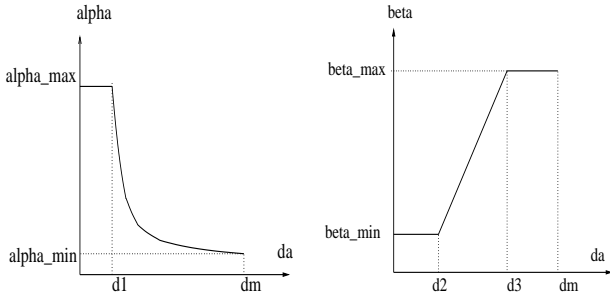


Fig. 1. α and β curves Vs d_a .

$$\beta = f_2(d_a) = \begin{cases} \beta_{min} & \text{if } d_a \leq d_2 \\ \kappa_3 + \kappa_4 d_a & \text{if } d_2 < d_a < d_3 \\ \beta_{max} & \text{otherwise.} \end{cases} \quad (3)$$

We let $f_1(\cdot)$ and $f_2(\cdot)$ be continuous functions and thus $\frac{\kappa_1}{\kappa_2 + d_1} = \alpha_{max}$, $\beta_{min} = \kappa_3 + \kappa_4 d_2$ and $\beta_{max} = \kappa_3 + \kappa_4 d_3$. Suppose d_m is the maximum average queueing delay and let $\alpha_{min} = f_1(d_m)$; then we also have $\frac{\kappa_1}{\kappa_2 + d_m} = \alpha_{min}$. From these conditions, we have

$$\begin{aligned} \kappa_1 &= \frac{(d_m - d_1)\alpha_{min}\alpha_{max}}{\alpha_{max} - \alpha_{min}} & \text{and} & \quad \kappa_2 = \frac{(d_m - d_1)\alpha_{min}}{\alpha_{max} - \alpha_{min}} - d_1, \\ \kappa_3 &= \frac{\beta_{min}d_3 - \beta_{max}d_2}{d_3 - d_2} & \text{and} & \quad \kappa_4 = \frac{\beta_{max} - \beta_{min}}{d_3 - d_2}. \end{aligned} \quad (4)$$

This specific choice is shown in Fig. 1.

We now describe the TCP-Illinois protocol in more details:

- All the features of TCP-NewReno except the AIMD algorithm are retained.
- In the congestion avoidance phase, the sender measures RTT T for each acknowledgement, and average the RTT measurements over the last W acknowledgements (one RTT interval) to derive the average RTT T_a . The sender records the maximum and minimum (average) RTT⁴ ever seen as T_{max} and T_{min} , respectively, and computes the maximum (average) queueing delay $d_m = T_{max} - T_{min}$ and the current average queueing delay $d_a = T_a - T_{min}$.
- The sender picks the following parameters: $0 < \alpha_{min} \leq 1 \leq \alpha_{max}$, $0 < \beta_{min} \leq \beta_{max} \leq 1/2$, $W_{thresh} > 0$, $0 \leq \eta_1 < 1$, $0 \leq \eta_2 \leq \eta_3 \leq 1$. The sender sets $d_i = \eta_i d_m$ ($i = 1, 2, 3$), computes κ_i ($i = 1, 2, 3, 4$) from (4), and computes α and β values from (2) and (3), respectively. The standard settings of these parameters are given in Section VI.
- $\alpha \leftarrow 1$ and $\beta \leftarrow 1/2$ if $W < W_{thresh}$.
- The κ_i ($i = 1, 2, 3, 4$) values are updated if T_{max} or T_{min} is updated. The α and β values are updated once per RTT.
- $W \leftarrow W + \alpha/W$ for each ACK.
- $W \leftarrow W - \beta W$, if in the last RTT there is packet loss detected through triple duplicate ACK.
- Once there is a timeout, the sender sets the slow start threshold to be $W/2$, enters slow start phase, and resets $\alpha = 1$ and $\beta = 1/2$, and α and β values are unchanged until one RTT after the slow start phase ends.

⁴The are two options here. In the default option, the maximum and minimum average RTTs are recorded. In an alternative option, the maximum and minimum instantaneous RTTs are recorded. These two options yield almost identical results, unless the delay signal is buried with noise and the noise is in a high level. Under this case, the default option is a better choice.

TCP-Illinois retains the fast recovery and fast retransmission features of NewReno in standard option. If the receivers support selective acknowledgement, TCP-Illinois can also back off its window size when packet loss is detected through selective ACK and adopt features from SACK TCP. However, the SACK support is not needed, since TCP-Illinois avoids heavy congestion effectively.

In addition to the above major features of the protocol, TCP-Illinois also contains another feature to improve the robustness against sudden fluctuations in delay measurements that can result from measurement noise, bursty packet arrival process, etc. To understand this feature, note that ideally once d_a becomes greater than d_1 , it should stay above d_1 until some users reduce their window sizes. However, due to bursty packet arrival process or measurement noise, it is possible for d_a to drop rather suddenly below d_1 before some users reduce their window sizes. In this case, we should not set $\alpha = \alpha_{max}$ unless we are really sure that the network is not in a congested state. Therefore, once $d_a > d_1$, we do not allow α to increase to α_{max} unless d_a stays below d_1 for a certain amount of time. TCP-Illinois chooses another parameter θ , and lets θ times RTT be this amount of time. The standard setting for θ is again given in Section VI.

We note that the adaptation of α is the key feature of TCP-Illinois, whereas the adaptation of β as a function of average queueing delay is only relevant in networks where there are non-congestion-related losses, such as wireless networks or extremely high speed networks. In wireless networks, some packet losses arise from channel fluctuations. In extremely high speed networks, congestion loss probability is so small that it is at the same level as or even smaller than the probability of packet transmission error at the link, and as a consequence, a non-trivial proportion of packet losses are from transmission error. For these non-congestion-related packet losses, we wish to avoid a sharp window size reduction. Then, the β adaptation of TCP-Illinois shows its advantage: although it still reduces window size, the reduction percentage is very small, since the queueing delay is very small.

IV. FAIRNESS AND STABILITY

In this section, we study the fairness and stability of TCP-Illinois. This involves both the intra-protocol fairness between different TCP-Illinois users and also inter-protocol fairness with standard TCP, i.e., the resource allocation between TCP-Illinois users and standard TCP users. We first develop a new stochastic matrix model for a class of general AIMD algorithms, which include standard TCP and TCP-Illinois as special cases, and then study the fairness and stability properties of these algorithms using this new model.

A. Stochastic Matrix Model for general AIMD Algorithms

There have been several recent papers on the stochastic matrix model of AIMD algorithms; see [2]–[4], [16], [25], [29]. We first provide an overview of this model, and then extend this model by modifying one of the assumptions in the earlier work. Throughout, we consider networks with a single

bottleneck link which uses Droptail, analyze the congestion avoidance phase only, and assume that all packet losses are caused by congestion.

Suppose a link with capacity C and queue limit B is shared by N users, indexed by i ($i = 1, 2, \dots, N$). User i has a transmission rate (or throughput) x_i , a window size W_i , a window increment parameter α_i , a window backoff factor β_i , and RTT T_i . We define $\mathbf{W} := [W_1, \dots, W_N]^T$, and $\mathbf{x} := [x_1, \dots, x_N]^T$. When the link is congested and one or more packets are dropped, we call this a congestion event, and denote by t_k the time at the k -th congestion event. At a congestion event, one or more flows see packet losses and backoff their window sizes, and we say that a loss event⁵ happens for these flows. For any variable v , we use $v(t)$ to denote its value at time t , use $v[k]$ (respectively, $v[k^+]$) to denote its value just before (respectively, after) the k -th congestion event, use $E[v]$ to denote the expected value of $v[k]$, and use \bar{v} to denote the average of all $v[k]$'s. Here, v could stand for W_i , x_i , α_i , β_i , T_i , \mathbf{W} , \mathbf{x} , as well as some other variables to be introduced later.

We now consider the congestion event k for Droptail queue. When congestion happens, the buffer is full, so every user experiences a maximum queueing delay $d_m = B/C$, and thus $T_i[k] \equiv \hat{T}_i := T_i^p + d_m, \forall k$, where T_i^p is the propagation delay of user i . At the congestion event, the instantaneous throughput for user i is $W_i[k]/\hat{T}_i$. Ignoring the burstiness of packet arrival process, we can assume that the outgoing packets from one particular user are evenly distributed along the path. For user i , altogether there are $W_i[k]$ packets, and thus the number of packets from user i queued in the link buffer should be $W_i[k]d_m/\hat{T}_i = x_i[k]d_m$. The sum of the queued packets from all users should be the link buffer limit B , and thus we have

$$B = \sum_{i=1}^N x_i[k]d_m = \sum_{i=1}^N x_i[k] \frac{B}{C}, \quad (5)$$

which leads to following equation:

$$\sum_{i=1}^N x_i[k] = C, \quad \forall k \in \{0, 1, 2, \dots\}. \quad (6)$$

As mentioned earlier, in this analysis we have ignored the burstiness of packet arrival process; if we had considered this burstiness, then $\sum_{i=1}^N x_i[k]$ would not be a constant, and would be either greater than or less than C . We now define $\Sigma = \{\mathbf{z} = [z_1, \dots, z_N]^T \in \mathbb{R}^N : z_i \geq 0, \sum_{i=1}^N z_i = C\}$; then Σ is the set of all possible $\mathbf{x}[k]$'s, and we call Σ the feasible set of $\mathbf{x}[k]$.

Between two consecutive congestion events, $W_i(t)$ is increased at rate $\alpha_i(t)/T_i(t)$, and thus

$$W_i[k+1] = W_i[k^+] + \int_{t_k}^{t_{k+1}} \frac{\alpha_i(t)}{T_i(t)} dt. \quad (7)$$

If we define

$$\tilde{T}_i[k] := \frac{\int_{t_k}^{t_{k+1}} \alpha(t) dt}{\int_{t_k}^{t_{k+1}} \frac{\alpha(t)}{T_i(t)} dt}, \quad (8)$$

⁵In our terminology, a congestion event is for a link, while a loss event is for an individual user.

then, we have

$$W_i[k+1] = W_i[k^+] + \frac{1}{\tilde{T}_i[k]} \int_{t_k}^{t_{k+1}} \alpha_i(t) dt. \quad (9)$$

For any user i and congestion event k , $\tilde{T}_i[k] \in [T_i^p, \hat{T}_i]$. In general, queueing delay is much smaller than propagation delay, and $\tilde{T}_i[k]$ varies in a very small range, so we can assume that $\tilde{T}_i[k] \equiv \tilde{T}_i, \forall k$.

At each congestion event k and for each flow i , we define the loss event random variable $D_i[k]$:

$$D_i[k] := \begin{cases} 1 & \text{if flow } i \text{ sees at least one packet loss,} \\ 0 & \text{otherwise,} \end{cases} \quad (10)$$

and define $\mathbf{D}[k] := [D_1[k], \dots, D_N[k]]^T$. Note that $D_i[k]$ and $D_j[k]$ are correlated, since $\sum_{i=1}^N D_i[k] \geq 1$.

With the loss event random variables defined, we have

$$W_i[k^+] = W_i[k](1 - \beta_i[k]D_i[k]). \quad (11)$$

Combining (9) and (11), and using the fact that $x_i[k] = W_i[k]/\hat{T}_i$, we have

$$x_i[k+1] = x_i[k](1 - \beta_i[k]D_i[k]) + \frac{1}{\tilde{T}_i \hat{T}_i} \int_{t_k}^{t_{k+1}} \alpha_i(t) dt. \quad (12)$$

Equations (6) and (12) describe the discrete-time stochastic model of all general AIMD algorithms.

B. Markov Chain for Identical $\alpha_i(t)$ and constant $\beta_i[k]$

We consider the class of AIMD algorithms which have the following properties: (i) $\alpha_i(t) = \alpha(t), \forall i$, where $\alpha(t)$ is the common window increment for all users at time t ; (ii) $\beta_i[k] \equiv \hat{\beta}, \forall i, k$, where $\hat{\beta}$ is a constant independent of i and k . This class includes standard TCP obviously, since it satisfies $\alpha_i(t) \equiv 1, \beta_i[k] \equiv \hat{\beta} = 1/2, \forall i, k$. This class also includes TCP-Illinois. First, $\alpha_i(t) = \alpha(t)$, since the queueing delay is the same for all users. Here, we ignore the differences in α among different flows due to feedback delays, since the queueing delays are averaged to compute α . Then, $\beta_i[k] \equiv \hat{\beta} = \beta_{max}, \forall i, k$, since the average queueing delay d_a is larger than the threshold parameter d_3 when congestion happens, if the parameters are carefully chosen. Recall our modeling of congestion events in the previous subsection: we know that the maximum queueing delay d_m is reached at each congestion event. Even considering the averaging process, $d_a[k]$ is close to d_m and still larger than d_3 .

For this class of AIMD algorithms, from (6), we have

$$\int_{t_k}^{t_{k+1}} \alpha(t) dt \sum_{i=1}^N (\tilde{T}_i \hat{T}_i)^{-1} = \sum_{i=1}^N \hat{\beta} D_i[k] x_i[k], \quad (13)$$

and thus

$$\int_{t_k}^{t_{k+1}} \alpha(t) dt = \frac{1}{\sum_{i=1}^N (\tilde{T}_i \hat{T}_i)^{-1}} \sum_{i=1}^N \hat{\beta} D_i[k] x_i[k]. \quad (14)$$

Define

$$\gamma_i := (\tilde{T}_i \hat{T}_i)^{-1} / \sum_{j=1}^N (\tilde{T}_j \hat{T}_j)^{-1}, \quad (15)$$

and $\gamma = [\gamma_1, \dots, \gamma_N]^T$. Then $\sum_{i=1}^N \gamma_i = 1$. We now have

$$x_i[k+1] = x_i[k](1 - \hat{\beta}D_i[k]) + \gamma_i \sum_{j=1}^N x_j[k] \hat{\beta}D_j[k]. \quad (16)$$

In vector form, we have

$$\mathbf{x}[k+1] = A[k]\mathbf{x}[k], \quad (17)$$

where

$$\begin{aligned} A[k] &= A(\mathbf{D}[k]) \\ &= \text{diag}(1 - \hat{\beta}D_1[k], \dots, 1 - \hat{\beta}D_N[k]) \\ &\quad + \gamma \hat{\beta}(D_1[k], \dots, D_N[k]). \end{aligned} \quad (18)$$

We see that $\mathbf{x}[k]$ forms a discrete-time Markov Chain on the continuous state space Σ . For any k , $A[k]$ is a non-negative, random, column stochastic matrix [5], [10]. The property of this Markov Chain is determined by the A matrix, and thus determined by $\mathbf{D}[k]$.

Note that, although $\alpha_i(t)$ determines the window curve, the recovery time after a congestion event, and the utilization of the bandwidth, once all users see the same $\alpha(t)$ at any time, $\alpha(t)$ does not influence the discrete-time Markov Chain at congestion events, and thus the exact form of $\alpha(t)$ is not important to understand the macroscopic fairness properties of this class of algorithms. So this stochastic matrix model applies to the entire class of such algorithms, and the special choice of $\alpha(t)$ in TCP-Illinois is not important when analyzing the fairness of TCP-Illinois. This special choice of $\alpha(t)$ indeed influences many other properties, such as efficiency and synchronization, as we will discuss later.

C. Stability and Fairness: General Case

In this subsection, we study the stability and fairness properties of the Markov Chain defined in (17) and (18). Let \mathbb{S} be the set of all non-empty subsets of $\{1, 2, \dots, N\}$, and suppose $s[k] \in \mathbb{S}$ is the set of users that experience a loss event at congestion event k . Define $\rho_s[k] := \text{Prob}(s[k] = s)$, where $s \in \mathbb{S}$. Then, the distribution of $D[k]$ is determined by the values of $\rho_s[k]$, $\forall s \in \mathbb{S}$. Let $q_i[k] := \text{Prob}(D_i[k] = 1)$, then $q_i[k] = \sum_{s: i \in s} \rho_s[k]$, and $q_i[k]$ denotes the window backoff probability of user i at congestion event k .

Most prior work assumed that $D_i[k]$ is independent of $\mathbf{x}[k]$, i.e., $\rho_s[k]$ is a constant independent of k , for each $s \in \mathbb{S}$, and $s[k]$ is identically independently distributed (i.i.d.). From this assumption, $q_i[k]$ is also a constant independent of k and $\mathbf{x}[k]$ for all user i . However, in reality, at different congestion events, a flow is more likely to see a loss event when it has a larger throughput than when it has a smaller throughput. Therefore, we modify the stochastic matrix model by allowing $\mathbf{D}[k]$ to be dependent on $\mathbf{x}[k]$, and allowing $q_i[k]$ and $\rho_s[k]$ to be functions of $\mathbf{x}[k]$ as well:

$$\rho_s[k] = \rho_s(\mathbf{x}[k]), \forall l, k, \text{ and } q_i[k] = q_i(\mathbf{x}[k]), \forall i, k. \quad (19)$$

We make the following assumption on ρ_s , q_i , and \mathbf{D} :

Assumption 1. (i) $\rho_s(\cdot)$ and $q_i(\cdot)$ are continuous functions in $\mathbf{x}[k]$. (ii) For any realization of the infinite length Markov

Chain defined in (17) and (18) and for any user i , $D_i[k] = 1$ for infinitely many k 's almost surely, i.e., for any $J > 0$,

$$\text{Prob}(D_i[k] = 0, \forall k \geq J) = 0, \forall i \in \{1, 2, \dots, N\}.$$

We now state the following theorem.

Theorem IV.1. Under Assumption 1, the Markov Chain defined in (17) and (18) has a unique invariant distribution, and starting from any initial state, the distribution of $\mathbf{x}[k]$ converges to this invariant distribution. Moreover, the Markov Chain is ergodic, i.e., for any continuous function $h(\cdot) : \Sigma \rightarrow \mathbb{R}$, $h(\bar{\mathbf{x}}[k])$, the time average of $h(\mathbf{x}[k])$, equals $E[h(\mathbf{x}[k])]$, the expected value of $h(\mathbf{x}[k])$ under the invariant distribution.

Proof. See [18], a longer version of this paper. \square

With the existence and uniqueness of the invariance distribution established, we now study the fairness among different users, i.e., the resource allocation under the invariant distribution. We have the following theorems on fairness.

Theorem IV.2. If N users sharing one link have homogeneous RTTs, under the unique invariant distribution of the Markov Chain defined in (17) and (18), all flows share the same expected throughput $E[x_i[k]]$.

Proof. When all users have the same RTT, \hat{T}_i , \tilde{T}_i and γ_i are the same for all users. Then, from (18), the A matrix does not depend on i or T_i . If we swap user i and user j , the Markov Chain is the same as if we do not swap user i and user j , but swap $x_i[0]$ and $x_j[0]$. From Theorem IV.1, we know that the invariant distribution is unique, independent of the initial condition. Therefore, the invariant distribution is unchanged if we swap user i and user j , and thus $E[x_i[k]] = E[x_j[k]]$, $\forall i, j \in \{1, 2, \dots, N\}$. \square

Theorem IV.3. If N users sharing one link have heterogeneous RTTs, under the unique invariant distribution of the Markov Chain defined in (17) and (18), the following equation holds:

$$\hat{T}_i \tilde{T}_i E[x_i[k] q_i(\mathbf{x}[k])] = E[W_i[k] \tilde{T}_i q_i(\mathbf{x}[k])] = C_1, \quad \forall i, \quad (20)$$

where C_1 is a constant independent of i .

Proof. Taking expectation of $x_i[k+1]$ given $\mathbf{x}[k]$ in (16), we have

$$E[x_i[k+1] | \mathbf{x}[k]] = x_i[k] - \hat{\beta} q_i(\mathbf{x}[k]) x_i[k] + \gamma_i \sum_{j=1}^N \hat{\beta} q_j(\mathbf{x}[k]) x_j[k]. \quad (21)$$

Under the invariant distribution,

$$\begin{aligned} E[x_i[k]] &= E[x_i[k+1]] = E[E[x_i[k+1] | \mathbf{x}[k]]] \\ &= E[x_i[k]] - \hat{\beta} E[q_i(\mathbf{x}[k]) x_i[k]] \\ &\quad + \gamma_i \hat{\beta} \sum_{j=1}^N E[q_j(\mathbf{x}[k]) x_j[k]]. \end{aligned} \quad (22)$$

So we have $E[x_i[k] q_i(\mathbf{x}[k])] / \gamma_i = \hat{T}_i \tilde{T}_i E[x_i[k] q_i(\mathbf{x}[k])] = E[W_i[k] \tilde{T}_i q_i(\mathbf{x}[k])]$ is independent of i and we have proved (20). \square

D. Models for $\rho_s(\cdot)$ and $q_i(\cdot)$

From Theorem IV.3, we see that the resource allocation depends on the form of $q_i(\cdot)$. If $q_i(\cdot)$ is constant, it is exactly the same as the prior results in [29]. In our model, $q_i[k]$ is allowed to be a function of $\mathbf{x}[k]$. We need to specify the $q_i(\cdot)$ function to further analyze the fairness property. Recall that the dependence of $q_i[k]$ on $\mathbf{x}[k]$ arises from the fact that a flow with a larger throughput is more likely to see a loss event than a flow with a smaller throughput. Accordingly, we make the following assumption:

Assumption 2. *At each congestion event, the total number of packets dropped is a random variable that takes values in $\{1, 2, \dots, M_{\max}\}$, and its distribution is independent of k . Furthermore, for any packet dropped at congestion event k , the probability that it belongs to flow i is $x_i[k]/C$.*

This assumption is justified by the following reasoning: since the total arrival rate is independent of k , so is the distribution for the total number of packets dropped; since at least one packet is dropped and only a finite number of packets are dropped, there are lower and upper bounds for the total number of packets dropped; since the probability of an arbitrary packet belonging to flow i is $x_i[k]/C$, so is the probability of a dropped packet belonging to flow i .

Lemma IV.1. *Assumption 1 holds given Assumption 2*

Proof. We first prove that Assumption 1 (i) holds. Let M be the random variable indicating the total number of packets dropped in one congestion event, let $P_M(m) = \text{Prob}(M = m)$ for all $m \in \{1, 2, \dots, \hat{M}\}$, and let $\hat{M} = E[M]$. Then, we have

$$\begin{aligned} q_i &= 1 - \text{Prob}(\text{no dropped packets from flow } i) \\ &= \sum_{m=1}^{M_{\max}} P_M(m) [1 - (1 - \frac{x_i[k]}{C})^m] = f(x_i[k]), \end{aligned} \quad (23)$$

where $f(x) := \sum_{m=1}^{M_{\max}} P_M(m) f_m(x)$, and $f_m(x) := 1 - (1 - \frac{x}{C})^m$. Both $f_m(x), \forall m$ and $f(x)$ are strictly increasing continuous functions in $x \in [0, C]$. Note that Assumption 1 allows $q_i[k]$ to be functions of all users' rates, while Assumption 2 further tells us that $q_i[k]$ is only a function of its own rate $x_i[k]$, and this relationship $f(\cdot)$ is common for all users.

We then study $\rho_s(\cdot)$. For a specific $s \in \mathbb{S}$, suppose $s = \{i_1, i_2, \dots, i_H\}$, where $1 \leq i_1 < i_2 < \dots < i_H \leq N$. Then, $\text{Prob}(s[k] = s | M = m) = 0$ if $m < H$. If $m \geq H$, we have

$$\begin{aligned} \rho_{s,m}(\mathbf{x}[k]) &:= \text{Prob}(s[k] = s | M = m) \\ &= \sum_{m_1, \dots, m_H} (\frac{x_{i_1}[k]}{C})^{m_1} \dots (\frac{x_{i_H}[k]}{C})^{m_H} \binom{m}{m_1, m_2, \dots, m_H}, \end{aligned}$$

where the summation is over all $m_h \geq 1, \forall h \in \{1, 2, \dots, H\}$, and $\sum_{h=1}^H m_h = m$. And we have

$$\rho_s(\mathbf{x}[k]) = \text{Prob}(s[k] = s) = \sum_{m=H}^{\infty} P_M(m) \rho_{s,m}(\mathbf{x}[k]). \quad (24)$$

So both $q_i(\mathbf{x}[k])$ and $\rho_s(\mathbf{x}[k])$ are continuous functions of $\mathbf{x}[k]$.

Next we prove Assumption 1 (ii). At each congestion event, at least one user will decrease its window size by at least 1, and thus $\sum_{i=1}^N x_i[k] \hat{\beta} D_i[k] \geq 1/T_m$, where $T_m =$

$\max_i \hat{T}_i$. Hence, if user i does not back off at congestion event k , $x_i[k+1] \geq \gamma_i/T_m$. Since at least one packet is lost at congestion event $k+1$ and the probability of a lost packet belonging to flow i is $x_i[k+1]/C$, we know that $q_i[k+1] \geq x_i[k+1]/C$. Thus $q_i[k+1] \geq \varepsilon := \gamma_i/(CT_m) > 0$. So the probability that user i backs off at least once in any two consecutive congestion events is lower bounded by $\varepsilon > 0$. As a consequence, for any user $i \in \{1, 2, \dots, N\}$, $D_i[k] = 1$ for infinitely many k 's almost surely. \square

Since Assumption 2 implies Assumption 1, we know that Theorem IV.1, Theorem IV.2 and Theorem IV.3 hold under Assumption 2 also. In the next subsection, we analyze the fairness property for the specific $f(x_i)$ function given by (23).

E. Synchronization and Fairness

From Theorem IV.3 and equation (23), the $f(\cdot)$ function uniquely determines the backoff behavior and the fairness property. Different $f(\cdot)$ functions lead to different backoff behaviors: for example, if $f(\cdot) \equiv 1$, the backoffs are completely synchronized; otherwise, they are not. The exact form of $f(\cdot)$ depends on the distribution of M , and is thus unknown if $P_M(\cdot)$ is unknown. However, we can bound $f(x)$ in general and approximate $f(x)$ for some special cases. Since

$$1 - \frac{x}{C} \geq (1 - \frac{x}{C})^m \geq 1 - \frac{mx}{C}, \quad \forall 0 \leq x \leq C,$$

we have

$$\frac{x}{C} \leq f(x) \leq \sum_{m=1}^{M_{\max}} P_M(m) [1 - (1 - \frac{mx_i[k]}{C})] = \frac{\hat{M}x}{C}. \quad (25)$$

Note that the loss event for a flow is the union of the events that each dropped packet belongs to this flow, therefore the bound on $f(x)$ in (25) is just the union bound. Since $x \geq 0$ and $f(x) \geq 0$ always, we have

$$\frac{x^2}{C} \leq E[xf(x)] \leq \frac{\hat{M}x^2}{C}.$$

From Theorem IV.3, we have

$$\frac{\hat{T}_j}{\hat{M}\hat{T}_j} \leq \frac{\bar{W}_i^2}{\bar{W}_j^2} \leq \frac{\hat{M}\hat{T}_i}{\hat{T}_i}, \quad \forall i \neq j. \quad (26)$$

The bounds are tight and $\bar{W}_i^2 \approx \bar{W}_j^2$ if \hat{M} is close to 1 (very light congestion). If the variance of $W_i[k]$ is much smaller than $(E[W_i[k]])^2$, then $\bar{W}_i^2 \approx (\bar{W}_i)^2$, and thus the average window sizes of all flows are almost the same under very light congestion case. From simulations which will be presented later, we observe that $\text{Var}(W_i[k]) \ll (E[W_i[k]])^2$ indeed. If \hat{M} is very large (heavy congestion), the bounds are meaningless and \bar{W}_i^2 and \bar{W}_j^2 can be significantly different.

We then consider the approximation of $f(x)$ under some special cases. When x is small, such that $mx/C \ll 1$, the probability that more than one dropped packets belong to one flow (with rate x) is very small, and thus the union upper bound is nearly reached: $f_m(x) \approx mx/C$. If $M_{\max}x/C \ll 1$, then $f(x) \approx \hat{M}x/C$, and $f(x) \propto x$. When x is large, such that $mx/c \gg 1$, then $f_m(x) \approx 1$. If $M \gg c/x$ with very high probability, then $f(x) \approx 1$. In general, $f(x)$ is a concave

curve, and $f(x) \propto x^\lambda$, where $0 \leq \lambda \leq 1$, and $\lambda \rightarrow 1$ as $M_{max} \rightarrow 1$, and $\lambda \rightarrow 0$ as $\hat{M} \rightarrow \infty$.

As we have mentioned, simulations show that the standard deviation of $x_i[k]$ is very small compared with $E[x_i[k]]$, so with high probability, $x_i[k]$ lies not far away from $E[x_i[k]]$. If the RTTs of different users do not differ significantly, $E[x_i[k]]$ is not significantly different from C/N , and the probability of $x_i[k] \ll C/N$ or $x_i[k] \gg C/N$ is very small. So if $M_{max} \ll N$, which we call ‘‘light congestion’’, almost always $Mx_i[k]/c \ll 1$, and thus $f(x_i[k]) \propto x_i[k]$ ($\lambda \approx 1$), and the window backoffs of different users are completely unsynchronized. If $\hat{M} \gg N$, which we call ‘‘heavy congestion’’, almost always $Mx/c \gg 1$, and thus $f(x) \propto 1$ ($\lambda \approx 0$), and the window backoffs are completely synchronized. In the middle of these two extreme cases, $0 < \lambda < 1$, and the window backoffs are partially-unsynchronized.

Plugging $f(x) \propto x^\lambda$ into (20), we get $\hat{T}_i \tilde{T}_i E[(x_i[k])^\lambda]$ is the same for all users. If the variance of $x_i[k]$ is small, and if the difference between \tilde{T}_i and \hat{T}_i is small also, we get the following fairness property:

$$\bar{x}_i \propto \frac{1}{\hat{T}_i^{1+\mu}}, \text{ and } \bar{W}_i \propto \frac{1}{\tilde{T}_i^\mu}, \quad (27)$$

where $\mu = (1 - \lambda)/(1 + \lambda)$. We know that $0 \leq \mu \leq 1$ in general; $\mu \approx 0$ for light congestion and completely unsynchronized window backoff; and $\mu \approx 1$ for heavy congestion and completely synchronized window backoff. So when the variance of $x_i[k]$ is small and when the RTTs are not significantly different, the fairness depends on the synchronization, which further depends on the heaviness of congestion. Light congestion leads to unsynchronized window backoff and equality of window sizes ($\mu = 0$); heavy congestion leads to synchronized window backoff and an inverse proportional relationship between window size and RTT ($\mu = 1$); and the general case ($0 < \mu < 1$) lies in the middle of these two extreme cases.

We then explore the factors that influence the distribution of M . Consider the homogeneous RTT case and suppose the system is slotted with each slot being one RTT. Since the pipe can hold at most $CT + B$ packets and W_i increases by α_i within each slot, $\sum_{i=1}^N W_i \in [CT + B + 1 - \sum_{i=1}^N \alpha_i, CT + B]$ in the slot just before congestion, and $\sum_{i=1}^N W_i \in [CT + B + 1, CT + B + \sum_{i=1}^N \alpha_i]$ in the slot of congestion. As a result, anywhere from 1 to $\sum_{i=1}^N \alpha_i$ packets could be dropped at one congestion event, and we know that the congestion is heavier if the increment before congestion is larger. Approximately, we can assume that M takes values from 1 to $\sum_{i=1}^N \alpha_i$ with equal probability, and thus $M_{max} = \max(1, \sum_{i=1}^N \alpha_i)$, and $\hat{M} = \max(1, (1 + \sum_{i=1}^N \alpha_i)/2)$. Since TCP-Illinois chooses very small $\alpha \ll 1$ just before congestion, $M_{max} \ll N$ and light congestion condition is satisfied. From this analysis, one advantage of TCP-Illinois is that it avoids heavy congestion and synchronized backoff, and it reaches a fair resource allocation between different users. On the contrary, convex curve algorithms, like HS-TCP, yield heavy congestion regularly, and this further causes synchronization and unfairness, as shown in Section VI.

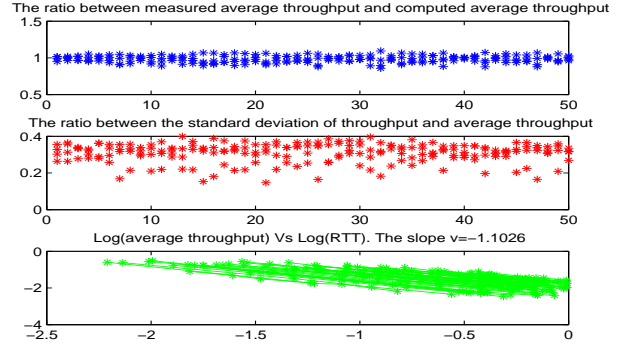


Fig. 2. $N = 4$, light congestion, $M \in [1, N/2]$, uniformly distributed.

We finally perform Matlab simulations to support our assumption of small $x_i[k]$ variance and validate our analysis on the relationship between heaviness of congestion and fairness. We have performed a large number of simulations on the evolution of the Markov Chain defined in (17) and (18). We vary N (the number of users) from 4 to 10. For each N , we select three probability distributions of M : (i) light congestion, M is uniformly distributed in $[1, N/2]$; (ii) medium level congestion, M is uniformly distributed in $[1, N]$; (iii) heavy congestion, M is uniformly distributed in $[1, 2N]$. For each scenario, we perform 50 simulations. For each simulation, $\mathbf{x}[0]$ and $T_i, \forall i$ are randomly generated initially, $\gamma_i, \forall i$ are computed by (15), and $M[k]$ and $s[k]$ are randomly generated according to Assumption 2 at each congestion event k , and thus each $A[k]$ and the sample path of the Markov Chain are derived. For each sample path, we average 1000 congestion events after the distribution converges, to compute average throughput \bar{x}_i , standard deviation of throughput $Std(x_i) := \sqrt{Var(x_i)}$, for each user i . We also compute x_i^* for all i by assuming that (27) holds for $\mu = 0$, i.e., all users share the same window size. We plot \bar{x}_i/x_i^* and $Std(x_i)/\bar{x}_i$ for all users and all simulations performed, and plot $\log(x_i)$ Vs $\log(T_i)$ for all users in each simulation. The results are shown from Fig. 2-7⁶. From the figures, we have the following observations: (i) the \bar{x}_i/x_i^* ratio is very close to 1 for light congestion, which indicates that all users share almost the same window size under light congestion, and as the congestion becomes heavy, the range of this ratio becomes wider and thus the difference between W_i 's becomes larger; (ii) the $Std(x_i)/\bar{x}_i$ ratio is always much smaller than 1 for any N and any heaviness of congestion, which supports our assumption of small variance; (iii) $\log(x_i)$ is linear with $\log(T_i)$, which validates the fairness property in (27): $x_i \propto 1/T_i^{1+\mu}$, where $\mu \in [0, 1]$, and μ increases as the congestion becomes heavy.

Remark 1. In equation (27), the average of W and x is over their values at the congestion events, and not over all time. Since a general AIMD algorithm can yield any window size curve, it is a challenging problem to compute

⁶Due to space limitation, we only provide the case for $N = 4$ and 10. For other values of N , the results turn out to be similar.

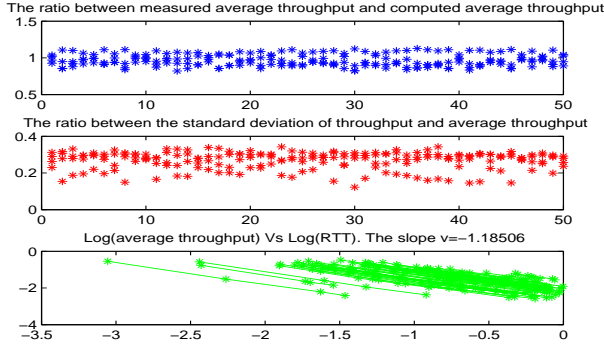


Fig. 3. $N = 4$, medium level congestion, $M \in [1, N]$, uniformly distributed.

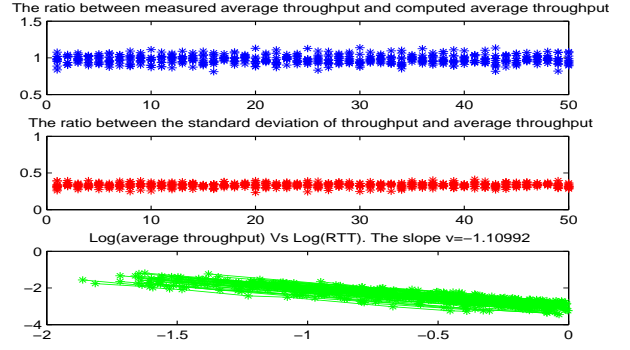


Fig. 5. $N = 10$, light congestion, $M \in [1, N/2]$, uniformly distributed.

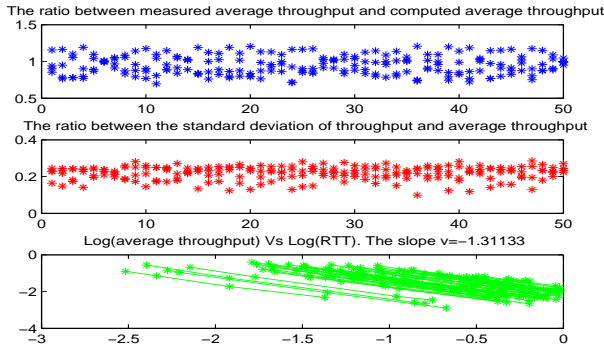


Fig. 4. $N = 4$, heavy congestion, $M \in [1, 2N]$, uniformly distributed.

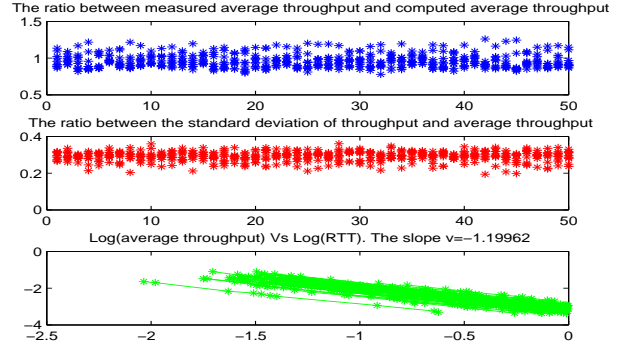


Fig. 6. $N = 10$, medium level congestion, $M \in [1, N]$, uniformly distributed.

the average W and x over all time, and it is an open problem whether the above conclusion on fairness holds for time averages of W and x . However, for general AIMD algorithms, since the all time average of W_i lies between $E[W_i[k]]$ and $E[W_i[k^+]] = E[(1 - \hat{\beta}_i[k])W_i[k]] \geq (1 - \beta_{max})E[W_i[k]]$, we know that time average of W for one user is similar to the average at congestion events. For TCP-Illinois, in particular, since x increases to near full utilization very quickly and stays around full utilization for a long time, the time average of x is very close to $E[x[k]]$, and thus the fairness property should hold approximately for time average also. From our ns-2 simulations in Section VI, we observe that the time averages of W for different users are also approximately the same.

F. Compatibility with the Standard TCP

We now consider the equilibrium allocation when TCP-Illinois coexists with TCP-Reno (NewReno and SACK are the same) and all flows share the same RTT⁷. If we use \hat{T} to denote the common maximum RTT for all flows, then $W_i[k] = x_i[k]\hat{T}$. The users are divided into two classes, Illinois user set \mathbb{I} and Reno user set \mathbb{R} , and $\alpha_i(t) = \alpha_{iL}(t)$, $\beta_i[k] = \beta_{iL}[k]$, $\forall t, \forall k, \forall i \in \mathbb{I}$, $\alpha_j(t) \equiv 1$, $\beta_j[k] \equiv 1/2$, $\forall t, \forall k, \forall j \in \mathbb{R}$.

⁷Due to lack of space, we considered here only the homogeneous RTT case. For the heterogenous RTT case, similar method can be used to derive the equilibrium allocation.

Define

$$\bar{\alpha}_i[k] = \frac{\int_{t_k}^{t_{k+1}} \frac{\alpha_i(t)}{T(t)} dt}{\int_{t_k}^{t_{k+1}} \frac{1}{T(t)} dt}. \quad (28)$$

Then, we know that $\bar{\alpha}_i[k] \equiv 1, \forall k, \forall i \in \mathbb{R}$ and $\alpha_{min} \leq \bar{\alpha}_i[k] = \alpha_j[k] \leq \alpha_{max}, \forall k, \forall i, j \in \mathbb{I}$, and we have

$$W_i[k+1] = (1 - \beta_i[k]D_i[k])W_i[k] + \bar{\alpha}_i[k] \int_{t_k}^{t_{k+1}} \frac{1}{T(t)} dt$$

From similar steps in (13) to (16), We have

$$W_i[k+1] = (1 - \beta_i[k]D_i[k])W_i[k] + \frac{\bar{\alpha}_i[k]}{\sum_{i=1}^N \bar{\alpha}_i[k]} \sum_{i=1}^N \beta_i[k]D_i[k]W_i[k]. \quad (29)$$

Suppose $\bar{\alpha}_i[k]$ is independent of $W_i[k]$, and define $\hat{\alpha}_i = E[\bar{\alpha}_i[k]]$, and $\hat{\beta}_i = E[\beta_i[k]]$. Then, $\hat{\alpha}_i = 1, \hat{\beta}_i = 1/2, \forall i \in \mathbb{R}$, and $\hat{\alpha}_j = \hat{\alpha}_k, \hat{\beta}_j = \hat{\beta}_k, \forall j, k \in \mathbb{I}$. We use α_{iL}^* and β_{iL}^* to denote the common $\hat{\alpha}_j$ and $\hat{\beta}_j$ values for TCP-Illinois users. If all packets dropped are due to congestion, and Droptail is used, $\beta_{iL}^* \approx \beta_{max}$. Taking conditional expectation of $W_i[k+1]$ given $W_i[k]$, we get

$$\begin{aligned} (\sum_{j=1}^N \hat{\alpha}_j)E[W_i[k+1]|W_i[k]] &= (\sum_{j=1}^N \hat{\alpha}_j)(W_i[k] \\ &- \hat{\beta}_i f(\frac{W_i[k]}{\hat{T}})W_i[k]) + \hat{\alpha}_i \sum_{j=1}^N \hat{\beta}_j f(\frac{W_j[k]}{\hat{T}})W_j[k]. \end{aligned}$$

Equating $E[W_i[k+1]]$ and $E[W_i[k]]$, we have

$$E[W_i[k]f(\frac{W_i[k]}{\hat{T}})] \frac{\hat{\beta}_i}{\hat{\alpha}_i} = C_2, \quad (30)$$

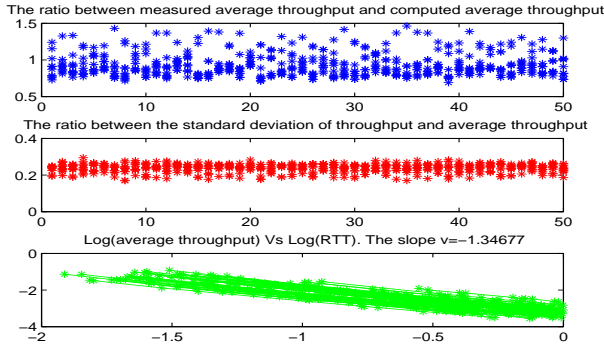


Fig. 7. $N = 10$, heavy congestion, $M \in [1, 2N]$, uniformly distributed.

where C_2 is a constant independent of user i . If the congestion is light and the backoff is unsynchronized ($\hat{M} \approx 1$ or $\hat{M} \ll C/x$), $f(x)$ is approximately proportional to x , and thus we have $\frac{\hat{\beta}_i}{\hat{\alpha}_i} E[(W_i[k])^2]$ is the same for all user i . And as $\text{Var}(W_i[k]) \ll (E[W_i[k]])^2$, approximately we have $\frac{\hat{\beta}_i}{\hat{\alpha}_i} (E[W_i[k]])^2$ is the same for all user i .

Since $\hat{\alpha}_i$ and $\hat{\beta}_i$ are the same within each protocol, we know that at equilibrium, all Reno users share the same average window size \bar{W}_R and all Illinois users share the same average window size \bar{W}_{IL} , and

$$\frac{\bar{W}_{IL}}{\bar{W}_R} \approx \sqrt{\frac{\alpha_{IL}^*}{2\beta_{IL}^*}} \approx \sqrt{\frac{\alpha_{IL}^*}{2\beta_{max}^*}},$$

In simulations presented later, for typical scenarios, $\sqrt{\alpha_{IL}^*}$ is slightly larger than 1 and β_{max} is usually picked to be $1/2$ for friendliness with TCP-Reno. Thus, \bar{W}_{IL} is usually slightly larger than \bar{W}_R . This means that in a network with both TCP-Illinois and TCP-Reno users, the TCP-Reno users will not suffer a significant degradation in performance. Furthermore, unlike TCP-Vegas which performs poorly when used with TCP-Reno, TCP-Illinois actually performs better than TCP-Reno, thus providing the right incentive for users to switch to TCP-Illinois.

V. TCP-ILLINOIS PROPERTIES

In Section II, we listed some requirements for the new TCP variant to satisfy, and in Section IV, we showed that TCP-Illinois maintains the intra protocol fairness the same way as standard TCP, satisfies the stability and scalability requirement, avoids heavy congestion, and is compatible with the current TCP. In this section, we consider the remaining requirements.

Since q_a increases with increasing W , α decreases with increasing W , and thus the W curve is concave. We can show that the curve is actually first linear, and then close to a parabola, and finally linear again. The proof is omitted due to space limitations and is available in [18]. In [18], we also show that TCP-Illinois achieves a better average throughput than standard TCP for any route buffer size B , and its average throughput increases as B increases, since compared with standard TCP, TCP-Illinois increases its rate to full

utilization faster and stays around full utilization longer (the length of time it stays around full utilization increases with increasing B). Thus the requirements of efficiency, router buffer independence, and incentive to switch are all met.

From Section IV, we see that convergence speed in k for TCP-Illinois is the same as that for standard TCP. So the response time is only determined by the time interval between two consecutive congestion events. We can show that this time interval of TCP-Illinois is similar to or smaller than that of standard TCP for a wide range of α_{min} values (see [18] for a proof), and thus the responsiveness requirement is also satisfied.

In lossy networks such as wireless networks, many packets are dropped not due to congestion. These packet drops greatly reduce the throughput for standard TCP, but for TCP-Illinois, the degradation is not as severe, since when a packet is dropped before congestion, the average queuing delay is always almost zero, and thus $\beta \approx \beta_{min}$ and $\alpha \approx \alpha_{max}$ always, and TCP-Illinois is essentially an AIMD algorithm with a larger $\alpha = \alpha_{max}$ and smaller β_{min} . Since $W \propto \sqrt{\alpha/\beta p}$, the ratio of the average window size of TCP-Illinois over that of standard TCP can be up to

$$\sqrt{\alpha_{max}/(2\beta_{min})}. \quad (31)$$

This improvement is significant. For example, if $\alpha_{max} = 9$, $\beta_{min} = 1/8$, then $W_{Illinois}$ can be up to $6W_{Reno}$.

VI. SIMULATION RESULTS

In this section, we provide *ns-2* simulation results to validate the properties of TCP-Illinois and compare its performance with TCP-Reno and HS-TCP. Throughout, one bottleneck link is shared by one or multiple users, which may choose TCP-Reno, HS-TCP, TCP-Illinois, or TCP-Vegas. For HS-TCP, all the default parameter settings are used. For TCP-Vegas, W increases if $diff < \gamma$ and decreases if $diff > \gamma$. For TCP-Illinois, without explicit explanation, we set $\alpha_{max} = 10$, $\alpha_{min} = 0.1$, $\beta_{max} = 1/2$, $\beta_{min} = 1/8$, $W_{thresh} = 10$, $\eta_1 = 0.01$, $\eta_2 = 0.1$, $\eta_3 = 0.8$, and $\theta = 5$.

A. Single User: Efficiency Property

We first perform simulations for a single user scenario, with $C = 100$ Mbps, $B = 100$ packets⁸, and $T_p = 100$ ms. The window sizes are plotted in Fig. 8. The simulations clearly demonstrate the concave nature of the curve of TCP-Illinois and show that TCP-Illinois achieves a larger average window size than TCP-Reno. For HS-TCP, we have chosen the Reno base, NewReno base and SACK base, and we have found that HS-TCP generates timeouts frequently for Reno and NewReno bases, and only works well if SACK is used. This supports our claim that HS-TCP causes heavy congestion. Numerically, the average send rates for TCP-Reno, SACK based HS-TCP⁹, and TCP-Illinois are 78.032, 87.324 and 91.304 Mbps, respectively. As a comparison to HS-TCP, if

⁸The packet size is 1000 bytes throughout.

⁹Henceforth, we mean SACK based HS-TCP when we mention HS-TCP without specifying its base.

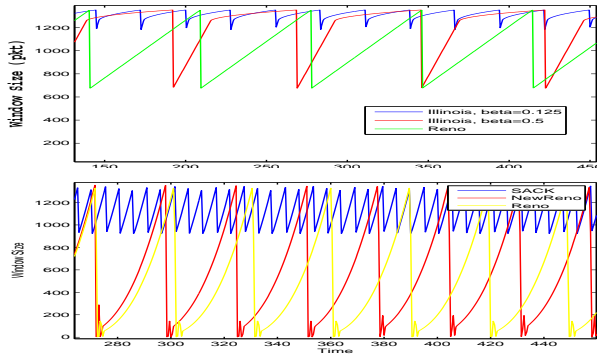


Fig. 8. Single user, TCP-Reno, HS-TCP, and TCP-Illinois. Top plot: Reno and Illinois ($\beta=0.125$ and $\beta=0.5$). Bottom plot: HS-TCP, with Reno, NewReno, and SACK bases.

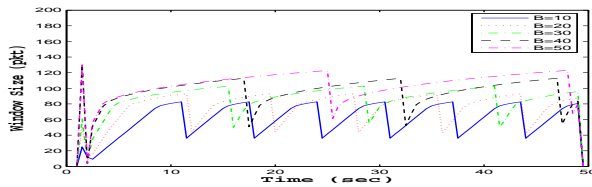


Fig. 9. The window size curve for different buffer sizes.

TCP-Illinois sets $\beta_{max} = \beta_{min} = 0.125$ ($\beta \approx 0.125$ for HS-TCP in this capacity range), then the average throughput is 95.727 Mbps.

We then study the effect of the buffer limit on the window size curve. We fix $C = 10$ Mbps, $T_p = 60$ ms, and vary B from 10 to 50 packets, and the window curve is plotted in Fig. 9. It is clear that as B increases, there is more time around full utilization and the average window size increases.

B. Multiple Users: Fairness Property

We now perform simulations for multiple ($N = 4$) users, which may choose Reno, HS-TCP or TCP-Illinois. We demonstrate the inter-protocol fairness (compatibility to Reno) of TCP-Illinois and HS-TCP in the homogeneous RTT scenario, where $C = 100$ Mbps, $B = 100$ packets, and $T_p = 100$ ms; and demonstrate the intra-protocol fairness of Reno, TCP-Illinois and HS-TCP in the heterogeneous RTT scenario, where C and B are unchanged, but the RTT s for the four flows are 60, 80, 100, and 120 ms, respectively. The average throughput in the homogeneous scenario and average window size in the heterogeneous scenario are plotted in Fig. 10. From this figure, we see clearly that TCP-Illinois is more fair to competing Reno user and large RTT users than HS-TCP.

To further demonstrate the performance of these protocols, we also plot the window curves of these simulations in Fig. 11 and Fig. 12.

C. Performance in Lossy/Wireless Networks

We then perform simulations for lossy/wireless links. It is a single link single user scenario, with the user choosing either TCP-Reno or TCP-Illinois, and the link randomly dropping packets with dropping probability p_d much larger

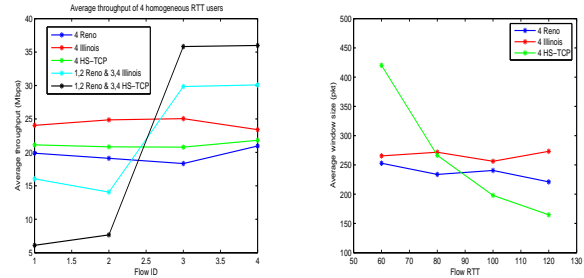


Fig. 10. Left: Average throughput of 4 homogenous RTT users. Right: Average window of 4 heterogeneous RTT users.

than the congestion loss probability (since p_d is large, the link is under utilized and there is no congestion loss at all in many cases). The capacity and buffer length of the link are 40 Mbps and 200 packets, respectively, and the propagation delay for the single user is 100 ms. Instead of choosing the default setting, the TCP-Illinois user sets $\eta_1 = 0.2$. We vary p_d values from 0.0005 to 0.05, and plot the average window size for TCP-Illinois and Reno and the ratio of these two multiplied by 20, as in the left plot of Fig. 13. From the plot, we see that $W_{Illinois} \approx 4W_{Reno}$ in most cases. From (31), the ratio should be $\sqrt{\alpha_{max}/(2\beta_{min})} = \sqrt{40} \approx 6.32$. The difference of the ratio between simulation and analysis can be explained if we observe the window curve plot of TCP-Illinois and Reno, as in the right plot of Fig. 13. From the window curve plot, we see that timeout happens frequently for TCP-Illinois user, since the increment amount α is very large before a packet loss happens. Equation (31) only considers the congestion avoidance phase, and timeout is the reason that $W_{Illinois}/W_{Reno}$ is around 4 instead of 6. Though, TCP-Illinois achieves a much better throughput than Reno in wireless networks.

D. Performance with Noisy RTT Measurement

We finally perform simulations to compare the performance of TCP-Illinois and TCP-Vegas when the delay measurement is inaccurate. We consider a single link and two user scenario. For the link, $C = 10$ Mbps and $B = 50$ packets (correspondingly, the maximum queueing delay $d_m = 40$ ms). For the users, either both choose TCP-Illinois or both choose TCP-Vegas, and the propagation delay is $T_p = 60$ ms for each user. We now suppose that there is an extra white noise term in the RTT measurement, denoted by n , and let n be uniformly distributed between $[0, 2\sigma]$ (the noise term is an extra delay due to reasons other than propagation and queueing, so it is nonnegative). Then, $RTT = T_p + d + n$, where d is the queueing delay. We vary the value of σ to vary the noise level and study the performance of TCP-Vegas and TCP-Illinois under noisy RTT measurement. For each protocol, we have two groups of simulations. In group one, both users face the noise term in the RTT measurement; and in group two, only one user faces the noise term and the other user measures RTT accurately. The average throughput of the users under different noise levels are plotted in Fig. 14. From Fig. 14, we see that as the noise level increases, TCP-

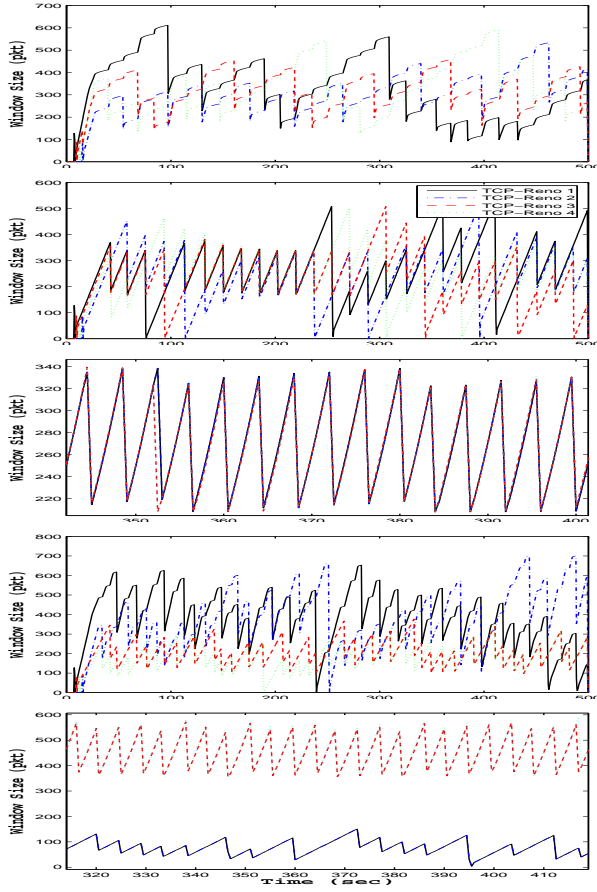


Fig. 11. Homogeneous RTT Users. First: 4 TCP-Illinois. Second: 4 Reno. Third: 4 HS-TCP. Fourth: 2 TCP-Reno and 2 TCP-Illinois. Fifth: 2 TCP-Reno and 2 HS-TCP. TCP-Illinois is demonstrated to avoid synchronization effectively and be compatible with TCP-Reno. HS-TCP is demonstrated to generate synchronization frequently and be unfair to TCP-Reno.

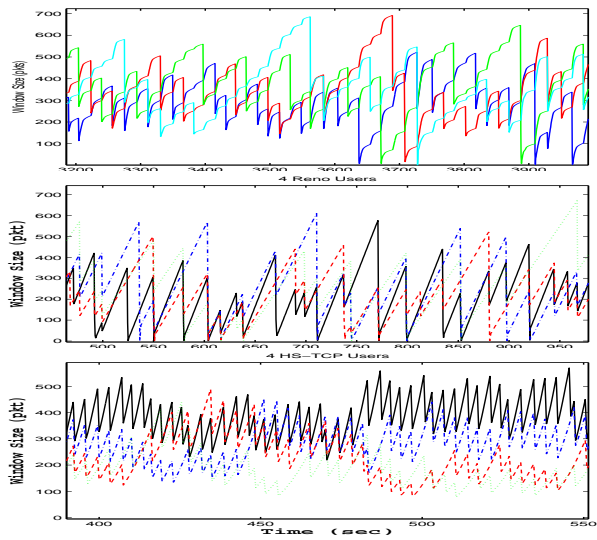


Fig. 12. Heterogenous RTT Users. Top plot: 4 TCP-Illinois. Middle plot: 4 Reno. Bottom plot: 4 HS-TCP. TCP-Illinois and Reno are fair to large RTT users, while HS-TCP is not.

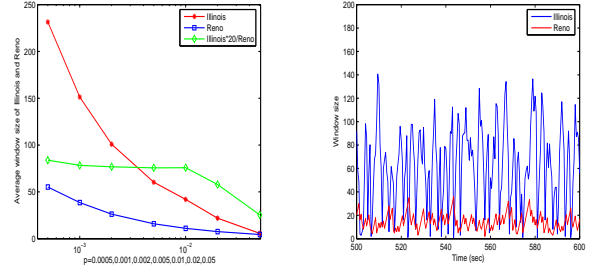


Fig. 13. Window sizes of TCP-Illinois Vs Reno in Lossy Networks. Left: Average window Vs p . Right: Window curve over t for $p = 0.005$.

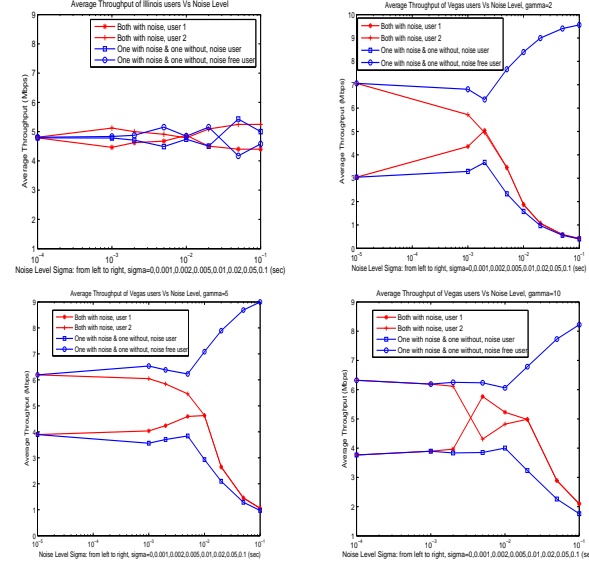


Fig. 14. Average Throughput Vs Noise Level. Top left: TCP-Illinois. Top right: Vegas, $\gamma = 2$. Bottom left: Vegas, $\gamma = 5$. Bottom right: Vegas, $\gamma = 10$.

Illinois is very robust to noise, while TCP-Vegas is not: there is a threshold on σ , which depends on the γ parameter. If the noise level exceeds this threshold, the performance is degraded significantly: if both users have noise terms, both get a much smaller throughput than the noise free case and the link is under utilized; if one user has the noise term and the other one does not, then the resource allocation is very unfair to the user with inaccurate RTT information. It is easy to explain the degradation when σ is larger than the threshold. At equilibrium, Vegas user satisfies

$$diff = \left(\frac{W}{T_p} - \frac{W}{T_p + d_a} \right) T_p = W \frac{d_a}{T_p + d_a} = \gamma \Rightarrow W^* = \gamma \frac{d_a^* + T_p}{d_a^*}, \quad (32)$$

where d_a^* and W^* are the equilibrium values of d_a and W . So W^* is an decreasing function of d_a^* , and d_a^* is a positive value such that the sum of $W^*/(T_p + d_a^*)$ over all users equals the capacity C . If $RTT = T_p + d_a + n$, then since n may hit zero, still we have $BaseRTT = T_p$, and (32) becomes to the following equation:

$$diff = \left(\frac{W}{T_p} - \frac{W}{T_p + d_a + n} \right) T_p = W \frac{d_a + n}{T_p + d_a + n} = \gamma \Rightarrow W = \gamma \frac{d_a + T_p + n}{d_a + n} \text{ and } \bar{W} \approx \gamma \frac{d_a + T_p + \bar{n}}{d_a + \bar{n}}, \quad (33)$$

where \bar{n} , \bar{d}_a and \bar{W} are the time average values of n , d_a and W . We see that if $\bar{n} = \sigma \leq d_a^*$, we can pick $\bar{d}_a = d_a^* - \sigma$ so that $\bar{W} = W^*$; and if $\bar{n} = \sigma > d_a^*$, then \bar{W} is definitely smaller than W^* . And when $\sigma > d_a^*$, as σ increases, \bar{W} decreases. So there exists a threshold of $\bar{n} = \sigma$ such that if the noise level is smaller than this threshold, the performance can be similar to the noise free case; and if the noise level is larger than this threshold, the performance is degraded and the degradation becomes more significant as σ increases. Since this threshold approximately equals to d_a^* and d_a^* is proportional to γ , we know that this threshold is also proportional to γ , as shown in Fig. 14.

VII. CONCLUSION

In this paper, we have considered some natural requirements for a new TCP protocol for high speed networks and have introduced a class of C-AIMD algorithms, which use loss to determine the *direction* and use delay to adjust the *pace* of window size change. This idea is rooted in the following two assumptions or understandings of the entire congestion control system: (i) delay is indeed a useful signal, i.e., congestion or packet loss is indeed correlated to delay information; (ii) delay is not an accurate signal, i.e., the correlation between loss and delay is weak. Combining these two, we should use loss as the primary signal and delay as the secondary signal. Using this idea, we have designed a specific protocol called TCP-Illinois, which achieves a concave window size curve and a better throughput than standard TCP, and maintains the fairness of standard TCP. Various properties of TCP-Illinois are studied, and TCP-Illinois is shown to satisfy all the requirements for an ideal high speed TCP variant.

To analyze the fairness property of TCP-Illinois, a new stochastic matrix model of general AIMD algorithms is introduced. Using this model, we have shown that TCP-Illinois leads to unsynchronized backoff and yields similar window size for different RTT users. There are still some open problems regarding the new model, however, which include: (i) the exact relationship between $E[W_i]$ and $E[W_j]$, $\forall i \neq j$ for the $N > 2$ heterogeneous users scenario; (ii), the relationship between $E[\bar{W}_i]$ and $E[\bar{W}_j]$, $\forall i \neq j$, where the expectation is over all time.

VIII. ACKNOWLEDGMENTS

We thank Professor F. Baccelli for helpful discussions.

REFERENCES

- [1] J. Ahn, P. Danzig, Z. Liu, and L. Yan. Experience with TCP vegas: Emulation and experiment. In *Proceedings of ACM SIGCOMM*, 1995.
- [2] F. Baccelli and D. Hong. Interaction of TCP flows as billiards. Technical Report, INRIA Rocquencourt, April 2002.
- [3] F. Baccelli and D. Hong. The AIMD model for TCP sessions sharing a common router. In *Proceedings of 39th Annual Allerton Conf. on Communication, Control and Computing*, October 2001.
- [4] F. Baccelli and D. Hong. AIMD, fairness and fractal scaling of TCP traffic. In *Proceedings of IEEE Infocom*, June 2002.
- [5] A. Berman and R. Plemmons. *Nonnegative Matrices in the Mathematical Sciences*. SIAM, 1979.
- [6] S. Biaz and N. Vaidya. Is the round-trip time correlated with the number of packets in flight. In *Proceedings Internet Measurement Conference (IMC)*, October 2003.

- [7] L. Brakmo, S. O'Malley, and L. Peterson. TCP vegas: New techniques for congestion detection and avoidance. In *Proceedings of ACM SIGCOMM Symposium*, pages 24–35, August 1994.
- [8] S. Floyd. Highspeed TCP for large congestion windows. Internet draft, draft-floyd-tcp-highspeed-01.txt, December 2003, Available at <http://www.icir.org/floyd/hstcp.html>.
- [9] S. Floyd and T. Henderson. The new reno modification to TCPs fast recovery algorithm. RFC 2582, 1999.
- [10] A. Graham. *Nonnegative Matrices and Applicable Topics in Linear Algebra*. Ellis Horwood Limited, Chichester, England, 1987.
- [11] V. Jacobson. Congestion avoidance and control. *ACM Computer Communication Review*, 18:314–329, August 1988.
- [12] V. Jacobson. Berkeley TCP evolution from 4.3-tahoe to 4.3-reno. In *Proceedings of the Eighteenth Internet Engineering Task Force*, July 1990.
- [13] C. Jin, D. Wei, S. H. Low, G. Buhrmaster, J. Bunn, D. H. Choe, R. L. A. Cottrell, J. C. Doyle, W. Feng, O. Martin, H. Newman, F. Paganini, S. Ravot, and S. Singh. FAST TCP: From theory to experiments, April 2003.
- [14] D. Katabi, M. Handley, and C. Rohrs. Congestion control for high bandwidth-delay product networks. In *Proceedings on ACM Sigcomm*, 2002.
- [15] T. Kelly. On engineering a stable and scalable TCP variant. Cambridge University Engineering Department Technical Report CUED/F-INFENG/TR.435, June 2002.
- [16] D. Leith and R. Shorten. Analysis and design of synchronised communication networks. *Automatica*, 41:725–730, 2005.
- [17] D. Leith, R. Shorten, and Y. Li. H-TCP: A framework for congestion control in high-speed and long-distance networks. HI Technical Report, August 2005. Available at <http://www.hamilton.ie/net/htcp/>.
- [18] S. Liu, T. Başar, and R. Srikant. TCP-illinois: A delay and loss-based congestion control algorithm for high-speed networks. Technical Report, UIUC, 2006. Available at http://www.ews.uiuc.edu/shaoliu/tcpillinois_full.pdf.
- [19] S. Liu, T. Başar, and R. Srikant. TCP-illinois: A loss and delay-based congestion control algorithm for high-speed networks.
- [20] J. Martin, A. A. Nilsson, and I. Rhee. Delay-based congestion avoidance for TCP. *IEEE/ACM Transactions on Networking*, pages 356–369, June 2003.
- [21] M. Mathis, J. Mahdavi, S. Floyd, and A. Romanow. TCP selective acknowledgement options. RFC 2018, April 1996, available at <http://www.icir.org/floyd/sacks.html>.
- [22] J. Mo, R. J. La, V. Anantharam, and J. C. Walrand. Analysis and comparison of TCP reno and vegas. In *Proceedings of IEEE INFOCOM*, 1999.
- [23] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose. Modeling TCP throughput: A simple model and its empirical validation. In *Proceedings of ACM SIGCOMM*, 1998.
- [24] R. S. Prasad, M. Jain, and C. Dovrolis. On the effectiveness of delay-based congestion avoidance. In *Proceedings of Second International Workshop on Protocols for Fast Long-Distance Networks*, 2004.
- [25] R. Shorten, F. Wirth, and D. Leith. A positive systems model of TCP-like congestion control: Asymptotic results. Submitted to *IEEE/ACM Transactions on Networking*, 2005.
- [26] W. Stevens. *TCP/IP Illustrated, Vol.1 The Protocols*. Addison-Wesley, 1994.
- [27] K. Tan, J. Song, Q. Zhang, and M. Sridharan. A compound TCP approach for high-speed and long distance networks. In *Proceedings of IEEE Infocom*, April 2006.
- [28] A. Tang, J. Wang, S. Low, and M. Chiang. Equilibrium of heterogeneous congestion control protocols. In *Proceedings of IEEE Infocom*, Miami, FL, March 2005.
- [29] F. Wirth, R. Stanojevic, R. Shorten, and D. Leith. Stochastic equilibria of AIMD communication networks. Accepted by SIAM J. Matrix Analysis and its Applications, 2005. Available at http://www.hamilton.ie/chris/SIMAX_july28.pdf.
- [30] L. Xu, K. Harfoush, and I. Rhee. Binary increase congestion control for fast long-distance networks. In *Proceedings of IEEE INFOCOM*, 2004.
- [31] A. Zanella, G. Procissi, M. Gerla, and M. Y. Sanadidi. TCP westwood: Analytic model and performance evaluation. In *Proceedings of IEEE Globecom*, 2001.