



NORTHWESTERN UNIVERSITY

Computer Science Department

Technical Report
NWU-CS-04-35
April, 2004

Modeling and Taming Parallel TCP on the Wide Area Network

Dong Lu Yi Qiao Peter A. Dinda Fabian E. Bustamante

Abstract

Parallel TCP flows are broadly used in the high performance distributed computing community to enhance network throughput, particularly for large data transfers. Previous research has studied the mechanism by which parallel TCP improves aggregate throughput and proposed a model to determine its upper bound when the network is not congested. In this work, we address how to predict parallel TCP throughput as a function of the number of flows without such constraints, as well as how to predict the corresponding impact on cross traffic. This combination allows us to answer the following question on behalf of a user: what number of parallel flows will give the highest throughput with less than a $p\%$ impact on cross traffic? We term this the maximum nondisruptive throughput. We begin by studying the behavior of parallel TCP in simulation to help derive a model for predicting parallel TCP throughput and its impact on cross traffic. Combining this model with some previous findings we derive a simple, yet effective, online advisor. We evaluate our advisor through simulation-based and wide-area experimentation.

Effort sponsored by the National Science Foundation under Grants ANI-0093221, ACI-0112891, ANI-0301108, EIA-0130869, and EIA-0224449. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation (NSF).

Keywords: Statistical TCP throughput characterization; TCP throughput prediction; TCP throughput monitoring; parallel TCP

Modeling and Taming Parallel TCP on the Wide Area Network

Dong Lu

Yi Qiao

Peter A. Dinda

Fabián E. Bustamante

{donglu,yqiao,pdinda,fabianb}@cs.northwestern.edu

Department of Computer Science, Northwestern University

Abstract—Parallel TCP flows are broadly used in the high performance distributed computing community to enhance network throughput, particularly for large data transfers. Previous research has studied the mechanism by which parallel TCP improves aggregate throughput and proposed a model to determine its upper bound when the network is not congested. In this work, we address how to predict parallel TCP throughput as a function of the number of flows without such constraints, as well as how to predict the corresponding impact on cross traffic. This combination allows us to answer the following question on behalf of a user: what number of parallel flows will give the highest throughput with less than a $p\%$ impact on cross traffic? We term this the maximum nondisruptive throughput. We begin by studying the behavior of parallel TCP in simulation to help derive a model for predicting parallel TCP throughput and its impact on cross traffic. Combining this model with some previous findings we derive a simple, yet effective, online advisor. We evaluate our advisor through simulation-based and wide-area experimentation.

I. INTRODUCTION

Data intensive computing applications require efficient management and transfer of terabytes of data over wide area networks. For example, the Large Hadron Collider (LHC) at the European physics center CERN is predicted to generate several petabytes of raw and derived data per year for approximately 15 years starting from 2005 [6]. Data grids aim to provide the essential infrastructure and services for these applications, and a reliable, high-speed data transfer service is a fundamental and critical component.

Recent research has demonstrated that the actual TCP throughput achieved by applications is, persistently, significantly smaller than the physical bandwidth “available” according to the end-to-end structural and load characteristics of the network [34], [23]. Here, we define *TCP throughput* as the ratio of effective data over its transfer time, also called *goodput* [30].

Parallel TCP flows have been widely used to increase throughput. For example, GridFTP [5], part of the Globus project [15], supports parallel data transfer and has been widely used in computational grids [6], [23].

A key challenge in using parallel TCP is determining the number of flows to use for a particular transfer. This number

affects both the throughput that the transfer will achieve and the impact that it will have on other traffic sharing links with these data flows. While there has been significant previous work on understanding and predicting parallel TCP performance, there is no analysis work or system that can support the following API call:

```
struct ParallelTCPChar {
    int    num_flows;
    double max_nondisruptive_thru;
    double cross_traffic_impact;
};

ParallelTCPChar *
TameParallelTCP(Address dest,
                 double maximpact);
```

Here, the user calls `TameParallelTCP()` with the destination of her transfer and the maximum percentage impact she is willing to have on cross traffic. The call evaluates the path and returns the number of parallel flows she should use to achieve the maximum possible throughput, while causing no more impact than the specified. We refer to this as the *maximum nondisruptive throughput (MNT)*.

The following sections address the implementation of such a function. With this in mind, we look for answers to the following questions:

- How does parallel TCP affect the throughput of the user’s transfer, the throughput of cross traffic, and the combined aggregate throughput, in different scenarios?
- How can these throughputs be predicted, online and with a small set of measurements, as functions of the number of parallel TCP flows?
- How can these predictions be used to implement the `TameParallelTCP()` function?

We begin by reviewing related work in Section II. In Section III we analyze parallel TCP throughput under different scenarios via simulations. We derive a prediction model for parallel TCP throughput and present results from an extensive Internet-based evaluation in Section IV. In Section V we outline a simple algorithm to estimate the effect of parallel TCP on cross traffic as a function of the number of flows. We evaluate our algorithm through simulation and later combine it with our approach to throughput prediction in order to implement the `TameParallelTCP()` call. Section VI presents our conclusions.

II. RELATED WORK

The available bandwidth of a path is defined as “the maximum rate that the path can provide to a flow, without reducing the rate of the rest of the traffic.” [18], [19]. Available bandwidth has been a central topic of research in packet networks over the years. To measure it accurately, quickly, and non-intrusively, researchers have developed a variety of algorithms and systems. Tools that measure either the bottleneck link capacity or the available bandwidth include cprobe [8], Remos [24], pathload [19], [20], NCS, and pipechar [21], among others [22], [10], [9], [36], [32], [18]. Most of these tools use packet pair or packet train techniques to conduct the measurements and typically take a long time to converge.

Previous research [22] has shown that, in most cases, the throughput that TCP achieves is considerably lower than the available bandwidth. Parallel TCP is one response to this observation. Sivakumar et al. [34] present P.Sockets, a library that stripes data over several sockets and evaluate its performance through wide-area experimentation. The authors concluded that this approach can enhance TCP throughput and, in certain situations, be more effective than tuning the TCP window size. Allcock et al. [6] evaluate the performance of parallel GridFTP data transfers on the wide-area, and applied GridFTP to the data management and transfer service in Grid environments.

Considerable effort has been spent on understanding the aggregate behavior of parallel TCP flows on wide area networks. Shenker et al [33] were first to point out that a small number of TCP connections with the same RTT and bottleneck can get their congestion window synchronized. Qiu et al. [30] studied the aggregate TCP throughput, goodput and loss probability on a bottleneck link via extensive ns2-based simulations. The authors found that a large number of TCP flows with the same round trip time (RTT) can also become synchronized on the bottleneck link when the average size of each TCP congestion window is larger than three packets. The reason for the synchronization is that, at the end of each epoch, the bottleneck buffer becomes full and each flow incurs a loss in the same RTT when it increments its congestion window. Since most flows have more than three outstanding packets before the loss, they can recover from the loss by fast retransmission and reduce the window by half, leading to global synchronization. Due to the global synchronization, all the flows share the resource fairly: in the steady state they experience the same loss rate, RTT and thus same bandwidth. Their findings are highly significant for our work, as they support our assumption that parallel TCP flows on the bottleneck link share the same loss rate.

The work most relevant to ours is that of Hacker et al [16]. The authors observe that parallel TCP increases aggregate throughput by recovering faster from a loss event when the network is not congested. The authors go on to propose a theoretical model for the upper bound of parallel TCP throughput for an uncongested path, i.e. the model is valid only if the network is not congested before and after adding in the parallel TCP flows. In contrast we attempt to predict throughput for both congested and uncongested paths, as well as to estimate

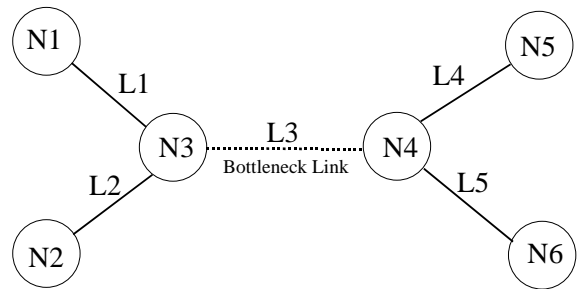


Fig. 1. Topology for simulations. Cross traffic goes from node N1 to N5, while parallel TCP flows go from node N2 to N6. Cross traffic and parallel TCPs share the same bottleneck link L_3 . Each simulation lasts 100 seconds with individual TCP cross traffic flows starting randomly during the first 8 seconds, and all parallel TCPs starting simultaneously at time 10 sec.

the impact of parallel TCP on cross traffic. Beyond this, our approach relies on only a small number of probes and no additional tools, while the model proposed in [16] is less practical because it requires loss rate measurements at each target number of parallel TCP flows, the cost of which is very high.¹

Hacker et al. [16] concluded that, in the absence of congestion, the use of parallel TCP flows is equivalent to using a large MSS on a single flow, with the added benefit of reducing the negative effects of random packet loss. They advise application developers not to use an arbitrary large number of parallel TCP flows, but conclude that it is difficult, if not impossible, to determine the point of congestion in the end-to-end path a priori, and therefore to decide on the proper number of parallel TCP flows. Our work attempts to advise the user as to how many flows to use given a limit on the impact to cross traffic, even when this cross traffic may be significant.

It is widely believed that, under congested situations, parallel TCP flows achieve better performance by effectively behaving unfairly, stealing bandwidth from cross traffic. This has prompted some researchers to propose modifying TCP in order to make it better suited for parallel transfers by considering both efficiency and fairness [11], [17]. We believe it will be difficult to persuade people to modify their TCP implementations just to use parallel TCP more fairly. By relying on our prediction tools, a user or administrator should be able to trade off a transfer’s throughput and its degree of impact on cross traffic, achieving what we refer to as the *maximum nondisruptive throughput (MNT)*. All without requiring modifications to pre-existing TCP implementations.

III. ANALYZING PARALLEL TCP THROUGHPUT

In this section, we use simulation to understand the behavior of parallel TCP under different scenarios. For all our simulation-based studies we make use of the ns2 network simulator [2].

A. Simulation Setup

In a simulation study on aggregate TCP throughput on a bottleneck link, Qiu et al. [30] developed a simple yet realistic

¹Conceivably, however, it may be possible to measure loss rates indirectly using software such as Web100 [26], [4].

topology model for wide-area Internet connections based on the Internet hierarchical routing structure (Figure 1). We adopt this same topology for our simulations. Each simulation is 100 seconds long, with cross traffic randomly starting during the first 8 seconds and parallel TCP flows all starting 10 seconds into the simulation. Cross traffic goes from N1 to N5, while parallel TCP flows go from N2 to N6. The bottleneck link is L3. We employ TCP Reno [12] for both cross traffic and parallel TCP flows, as this implementation makes use of the most widely deployed TCP congestion control algorithm.² Both DropTail and Random Early Detection (RED) [14] queue management policies are studied as they are the most commonly used queue management policies on the Internet. DropTail and RED have similar performance in most our simulations. The exception is in Scenario 1. Here, when there are more than 10 cross traffic flows, the cross traffic dominates the queue and starves the parallel TCP flows under DropTail policy. Unless otherwise noted, we show results for the DropTail policy.

We use TCP flows as cross traffic because of TCP’s dominance in the current Internet, as reported in the work by Smith et al. [35], work in which TCP accounted for 90-91% of the packets and about 90-96% of the bytes transferred in traces collected in 1999-2000 from a educational institution (UNC) and a research lab (NLANR).

We analyze Parallel TCP throughput under a variety of representative scenarios including a typical slow connection such as cable or DSL (Scenario 1), a coast-to-coast high-speed Internet connection (Scenario 2) and a current (Scenario 3) and next generation global-scale Internet connections (Scenario 4). Two additional scenarios (Scenarios 5 and 6) are used to represent cases where the TCP buffer has not been appropriately tuned [37]. Figure 2 summarizes the different simulation scenarios. For each scenario, we simulate from 1 to 31 parallel TCP flows with 5, 10, 15, 20, 25 and 30 random TCP cross traffic flows.

B. Simulation results

Figures 3 to 8 plot the aggregated throughput of parallel TCP as a function of the number of flows used for the different scenarios. Plots are shown both without (left graph) and with (right graph) cross traffic. In the latter case, we also plot the cross traffic’s and total throughput, i.e. the sum of both the parallel TCP and cross traffic throughputs.

Figure 3 shows our results for Scenario 1, used to represent a typical slow connection. We show five cross traffic flows in this case. It is clear from the graphs that, with such a low-latency/low-bandwidth connection, the primary benefit from parallel TCP comes from being able to steal bandwidth from the existing cross traffic.

The results for Scenario 2, representing a current coast-to-coast connection with low latency and medium bandwidth, are shown in Figure 4. As it can be seen from the plots, there are some limited benefits from using parallel TCP without competition in this scenario. In the presence of cross traffic, however, parallel TCP is an even stronger competitor. Notice

also how parallel TCP allows us to increase overall throughput, albeit marginally.

Figure 5 illustrates the benefits of parallel TCP in Scenario 3, a long latency, medium bandwidth link representing a current global-scale, fast Internet connection. In this case there are significant benefits to using parallel TCP even in the absence of cross traffic. The differences in the performance of parallel TCP under scenarios 2 and 3, without cross traffic, can be explained using Hacker’s theory [16]: parallel TCP recovers faster than single TCP when there is a time out. This effect is more important as the RTT increases, because the time out will be longer and single TCP cannot recover fast enough. As in our previous scenario, parallel TCP can aggressively steal bandwidth from the existing cross traffic, this time significantly increasing the overall throughput.

The benefits of using parallel TCP, with and without cross traffic, are very clear under Scenario 4 as Figure 6 shows. The additional throughput in the presence of cross traffic, is mainly due to the increase in overall throughput.

The advantage of parallel TCP is even more significant in the two scenarios representing mistuned TCP buffers. Figure 7 shows this advantage for Scenario 5, a high bandwidth and high latency link with a small socket buffer size. The benefits of parallel TCP are quite obvious, regardless of the amount of cross traffic. Furthermore, these gains come at no cost to the existing cross traffic. Parallel TCP gains performance not only by recovering faster after a time out, but also by providing an effectively larger buffer size. Note that the throughput of parallel TCP flows eventually flatten out as more flows are added into the simulation.

Similar benefits from parallel TCP can be observed in our last scenario (Figure 8). As with Scenario 5, parallel TCP can significantly improve throughput regardless of the degree of traffic. In this case, the impact on cross traffic increases with increasing parallelism, but remains relatively flat.

C. Observations

The dramatically different behaviors shown in the previous section clearly illustrate the challenges in providing a sound `TameParallelTCP()`-like call. The parallel TCP and cross traffic throughput curves adopt a wide range of forms, depending on the topology of the network and the configuration of endpoints. In addition, even if one were to disregard the almost prohibitively high costs of directly measuring these curves, the cross traffic impact would be very difficult to determine.

IV. MODELLING AND PREDICTING PARALLEL TCP THROUGHPUT

In this section we combined our simulation work together with our and others’ analytic treatment of TCP performance, to develop a model that can be used to predict the throughput of parallel TCP flows.

A. Algorithm

Mathis et al. [27] developed a simple model for single flow TCP Reno throughput on the assumption that TCP’s performance is determined by the congestion avoidance algorithm

²Comparable results were obtained using TCP Tahoe.

Scenario	L_3 latency	L_3 Bandwidth	L_1, L_2 Bandwidth	L_4, L_5 Bandwidth	TCP buffer
1	20 ms	1.5 Mbps	10 Mbps	10 Mbps	\geq Bandwidth*RTT
2	20 ms	100 Mbps	1000 Mbps	1000 Mbps	\geq Bandwidth*RTT
3	50 ms	100 Mbps	1000 Mbps	1000 Mbps	\geq Bandwidth*RTT
4	50 ms	1000 Mbps	10000 Mbps	10000 Mbps	\geq Bandwidth*RTT
5	50 ms	1000 Mbps	10000 Mbps	10000 Mbps	60 KB
6	20 ms	100 Mbps	1000 Mbps	1000 Mbps	60 KB

Fig. 2. Bandwidth and latency configuration for different scenarios. The latency for L_1 and L_2 is fixed at 4 milliseconds, while the latency for L_4 and L_5 is fixed at 5 milliseconds. The buffer size on each node is fixed at 25 packets. Both DropTail and RED queue management policies are simulated.

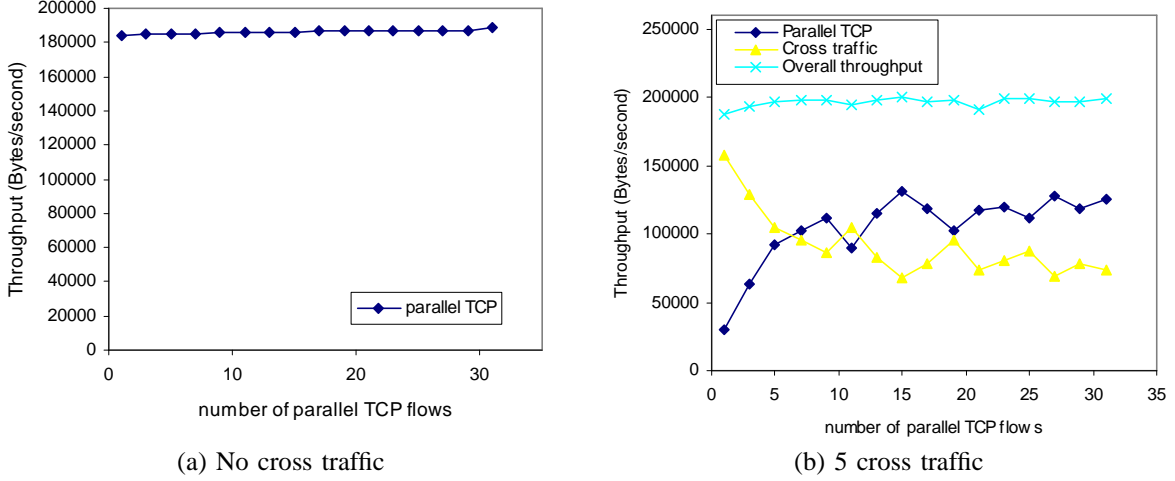


Fig. 3. Simulation results for scenario 1: latency of L_3 is 20 ms; bandwidth of L_3 is 1.5 Mbps; TCP buffer is properly tuned. Refer to Figure 2 for details.

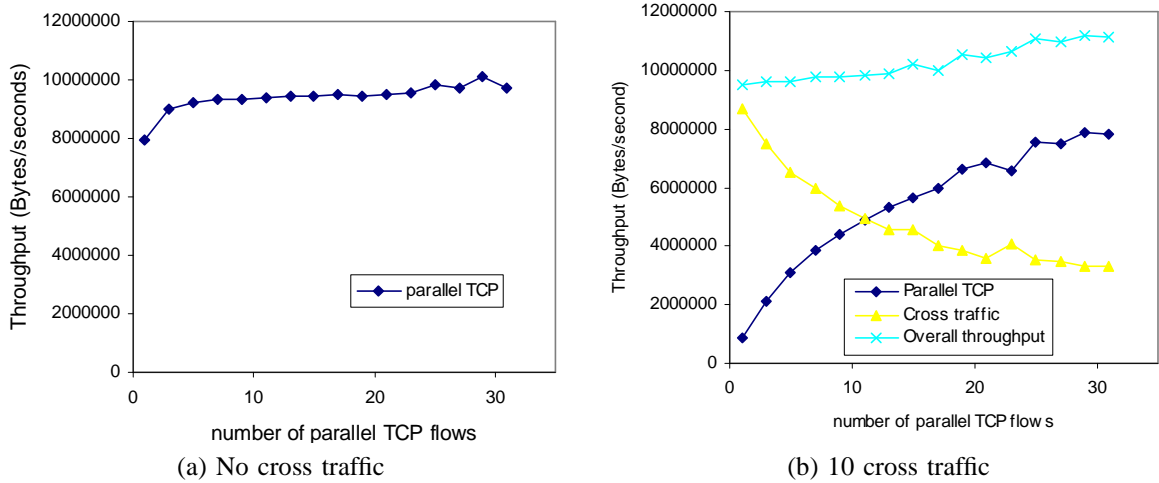


Fig. 4. Simulation results for scenario 2: latency of L_3 is 20 ms; bandwidth of L_3 is 100 Mbps; TCP buffer is properly tuned. Refer to Figure 2 for details.

and that retransmission timeouts are avoided:

$$BW = \frac{MSS}{RTT \sqrt{\frac{2bp}{3}}} \quad (1)$$

Here, p is the loss rate or loss probability, and b is the number of packets that are acknowledged by a received message. MSS and RTT are the maximum segment size and round trip time respectively.

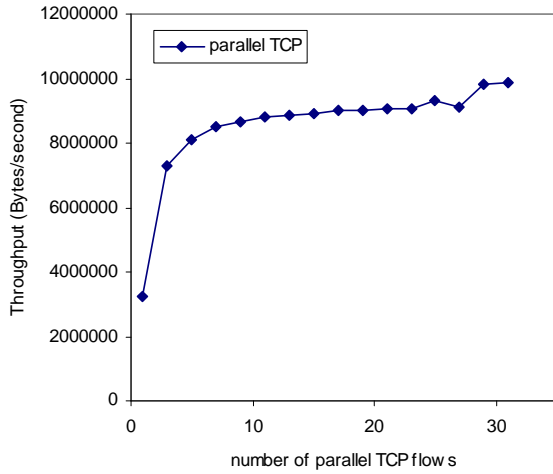
Padhye et al. [29] developed an improved single flow TCP Reno throughput model that considers timeout effects. Assuming that the TCP buffer is not the bottleneck (i.e., that the socket buffer size is large or “rightsized” [13]), their model

is

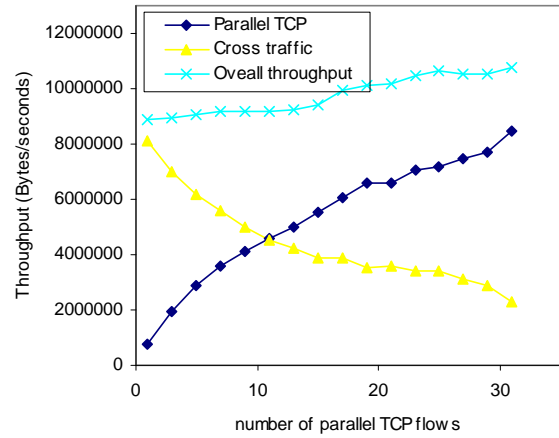
$$BW = \frac{MSS}{RTT \sqrt{\frac{2bp}{3}} + T_0 \min(1, 3\sqrt{\frac{3bp}{8}})p(1 + 32p^2)} \quad (2)$$

where T_0 is the timeout.

Bollinger et al [7] show that equations 1 and 2 are essentially equivalent with packet loss rates less than 1/100, something validated on the current Internet by Zhang et al [38]. Hacker et al. [16], based on Bollinger’s findings, present a model for the upper bound of the throughput of n parallel TCP flows. The authors assume that both MSS and RTT are stable. Hacker’s

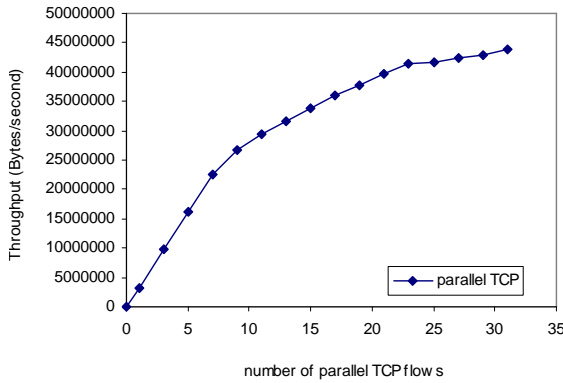


(a) No cross traffic

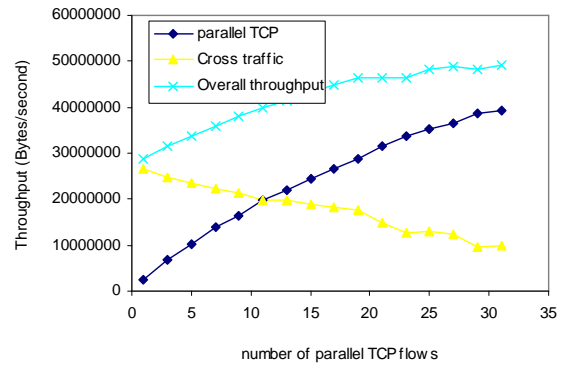


(b) 10 cross traffic

Fig. 5. Simulation results for scenario 3: latency of L_3 is 50 ms; bandwidth of L_3 is 100 Mbps; TCP buffer is properly tuned. Refer to Figure 2 for details.

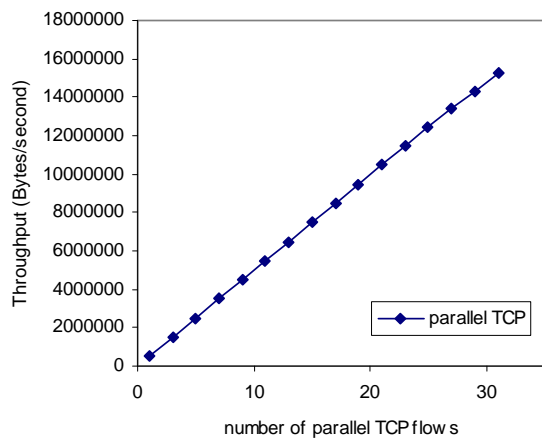


(a) No cross traffic

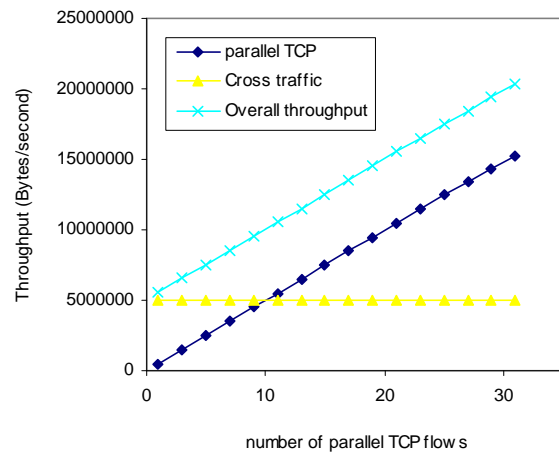


(b) 10 cross traffic

Fig. 6. Simulation results for scenario 4: latency of L_3 is 50 ms; bandwidth of L_3 is 1000 Mbps; TCP buffer is properly tuned. Refer to Figure 2 for details.



(a) No cross traffic



(b) 10 cross traffic

Fig. 7. Simulation results for scenario 5: latency of L_3 is 50 ms; bandwidth of L_3 is 1000 Mbps; TCP buffer is not properly tuned. Refer to Figure 2 for details.

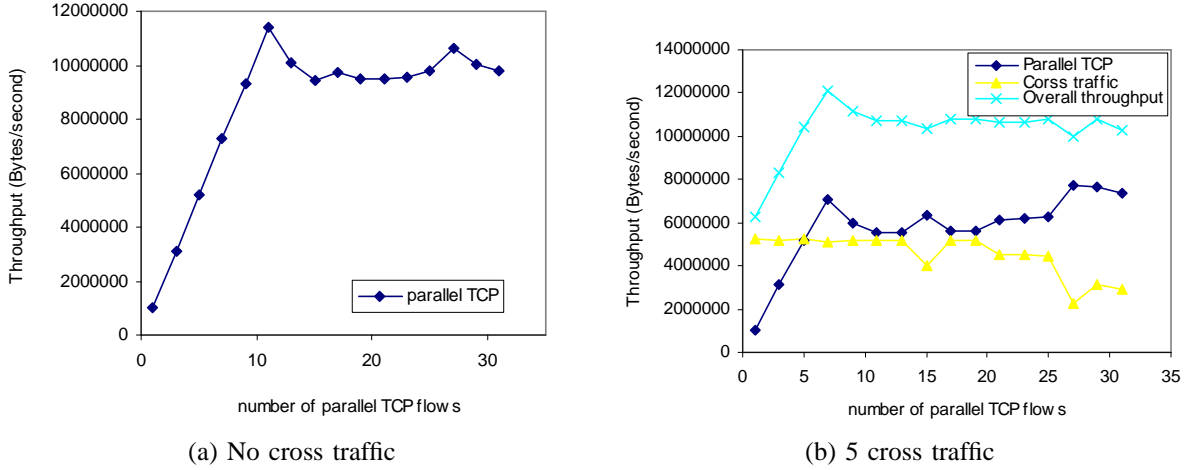


Fig. 8. Simulation results for scenario 6: latency of L_3 is 20 ms; bandwidth of L_3 is 100 Mbps; TCP buffer is not properly tuned. Refer to Figure 2 for details.

upper bound model can be stated as:

$$BW_n \leq \frac{MSS}{RTT} \sum_{i=1}^n \frac{1}{\sqrt{p_i}} \quad (3)$$

where p_i is the packet loss rate for flow i . This upper bound is useful when the network is not congested, but has limited utility otherwise. In addition, as already discussed in Section II, the online use of this model would require expensive measurements of different loss rates p_i at each level of parallelism.

In our model, we introduce the notion of the number of cross traffic flows, m , and assume that this does not change dramatically over significantly large time periods. Based on previous research [39] and our earlier work in characterizing, modelling, and predicting single flow TCP throughput [25], these appear to be reasonable assumptions. In addition, we assume that all of the parallel TCP flows see the same loss rate and have the same RTT , although both are functions of n and m . These two assumptions have been independently verified [30], as discussed in Section II. We denote with p_n the loss rate after adding n parallel TCP connections, and with RTT_n the round trip time at this point.

Along different paths, the value of MSS can vary ranging from the default 536 bytes to much larger values (for example to support 9000 byte Ethernet jumbo frames on LANs). Our prediction model does not depend on the a priori knowledge of MSS . We do assume, however, that this value does not change after connection establishment. This is a valid assumption as both sides with either use path MTU discovery at connection establishment time [28] or use the default 576 byte path MTU. MSS will directly follow from this.

Based on Equation 1 and the assumptions discussed above, we developed the following parallel TCP throughput model that essentially sums n TCP flows:

$$BW_n = \frac{MSS}{RTT_n} \frac{n}{\sqrt{p_n}} \frac{c_1}{\sqrt{\frac{2b}{3}}} \quad (4)$$

The TCP flows share the same RTT and loss rate and thus the same throughput. Both p_n and RTT_n are actually functions of

n and m . Given that we assume m is stable during a period of time, we treat them as functions of n alone. c_1 is a constant in the range $(0, 1]$ that we use to represent the effects of TCP timeouts. In the following, we assume that c_1 is stable for a path over at least short periods. Our model doesn't require the knowledge of c_1 .

If we had a model that could compute the relationship between p_n , RTT_n and the parallelism level n based on a small set of measurements, we could then use Equation 4 to predict the throughput for any parallelism level. This is in essence what we do. We developed several parametric models for this relationship based on measurements.

There is a trade-off between the sophistication of the model and the number of measurements needed to fit it. Recognizing this trade-off, we believe the best model for the relationship is a partial order two polynomial:

$$p_n \times RTT_n^2 = a \times n^2 + b \quad (5)$$

Here a and b are parameters to be fit based on measurements. In order to use the model in practice, we have to actively probe a path with two different levels of parallelism. The procedure is derived as follows.

Note that in Equation 4, $C = \frac{c_1}{\sqrt{\frac{2b}{3}}}$, and MSS are all constants under our assumptions. We define a new variable p'_n :

$$p'_n = p_n \frac{RTT_n^2}{C^2 MSS^2} = a'n^2 + b' \quad (6)$$

Combining Equations 4 and 6, we obtain:

$$BW_n = \frac{n}{\sqrt{p'_n}} \quad (7)$$

Based on Equation 7, we could use the TCP throughput at two different parallel levels to predict the TCP throughput at a third level. Let n_1 and n_2 be the two parallel levels that are probed:

$$BW_{n_1} = \frac{n_1}{\sqrt{p'_{n_1}}} = \frac{n_1}{\sqrt{a'n_1^2 + b'}} \quad (8)$$

and

$$BW_{n_2} = \frac{n_2}{\sqrt{p'_{n_2}}} = \frac{n_2}{\sqrt{a'n_2^2 + b'}} \quad (9)$$

From which we can determine:

$$a' = \frac{\frac{n_2^2}{BW_{n_2^2}} - \frac{n_1^2}{BW_{n_1^2}}}{n_2^2 - n_1^2} \quad (10)$$

and

$$b' = \frac{n_1^2}{BW_{n_1^2}} - a'n_1^2 \quad (11)$$

By substituting a' and b' in Equation 6 with the expression in Equation 10 and 11, we can now predict the TCP throughput for a third level of parallelism based on Equation 7.

Notice how our prediction requires only *two* TCP throughput probes, one for each of the two different parallel levels (n_1 and n_2). Both the probing and the calculation process are simple and incur little overhead, the majority of which lies in the communication cost of the two probes.

Equation 5 is an empirical approximation. In addition, we tested the effectiveness of a linear model,

$$p_n \times RTT_n^2 = a \times n + b \quad (12)$$

and a full order 2 polynomial that requires an additional (third) probe.

$$p_n \times RTT_n^2 = a \times n^2 + b \times n + c \quad (13)$$

We measured the performance of these three alternatives in a wide-area testbed [3], and found that (1) Equations 5 and 13 are better models than Equation 12, but (2) the full order two polynomial model is not significantly better than Equation 5 and can, indeed be sometime worse due to its sensitivity to sampling errors caused by small network fluctuations. Another problem with the full order two polynomial model is that it is sensitive to the choice of probe parallelism. Finally, it requires three probes instead of the two needed for the linear and partial polynomial models. As a result, we use Equation 5 for our system and the discussion in the rest of the paper, unless otherwise noted.

B. Evaluation

We evaluated our model extensively through online experimentation on PlanetLab [3], a planetary-scale testbed. We randomly choose 41 distinct end-to-end paths with end nodes located in North America, Asia, Europe and Australia. For each path, we conduct 10 rounds of experiments using Iperf [1] to obtain our measurements. A round of experiment starts with two probes for prediction purposes, immediately followed by parallel TCP transfers with up to 30 parallel TCP flows.

We adopt the *mean relative error* as our performance metric. Relative error is defined as:

$$relativeerror = \frac{prediction - measurement}{measurement} \quad (14)$$

Mean relative error on a path is the average of all the relative prediction errors on the path. Mean relative error for a given number of parallel TCP flows is the average of the relative prediction errors of all the experiments for that number of parallel TCP flows.

Figure 9 shows two examples of prediction using our model. The graphs show the actual and predicted throughput (based on measurements at $n_1 = 1$ and $n_2 = 10$). It can be seen

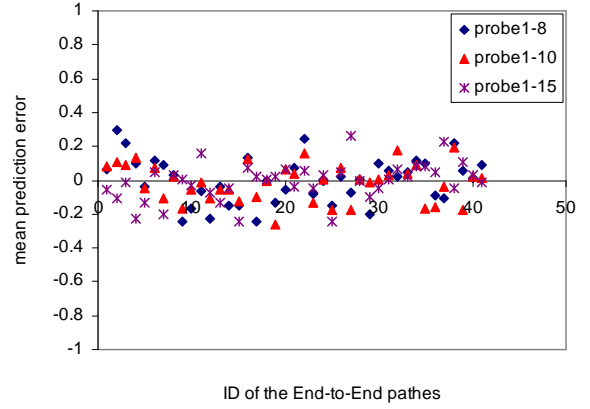


Fig. 11. Prediction sensitivity to the selection of probes.

that, for Example 1, predictions made based on the partial order 2 and full order 2 polynomials are virtually identical and have similar accuracy, while the prediction curve derived using the linear model deviates significantly from the measurement curve. In our second example, the prediction made using the partial order 2 polynomial and the linear model are virtually identical and equally accurate. The prediction curve generated by the full order 2 polynomial, however, deviates significantly from the measurement curve.

Figure 10 shows the performance of our parallel TCP throughput predictor using two probes at parallelism levels $n_1 = 1$ and $n_2 = 10$ for a wide range of PlanetLab pairs located all over the world. Only the partial order 2 polynomial model is used here. Each row in the table shows our relative prediction error for a particular Internet path between two hosts. The prediction quality is characterized by the mean and standard deviation of the relative errors at each of the different parallel levels (ranging between 1 and 30). The results presented in this table are quite encouraging: in most cases, our predictions guarantee us a small mean and standard deviation of relative prediction errors.

Our predictor is relatively insensitive to the particular level of parallelism for the probes. Figure 11 shows the mean relative error for our predictor using (1, 8), (1, 10) and (1, 15) parallel probes. We can see that we obtain similar performance in all cases. Of course, it is important not to use parallelism levels that are too close together (such as (1, 2)), as such probes are very sensitive to small fluctuations in the network or the existing cross traffic.

As it can be seen from Figure 12, the mean relative error for a given number of parallel TCP flows is not related to the number of parallel TCP flows. The figure, a scatter plot of the mean relative error versus the number of parallel TCP flows, shows no clear trend. The correlation coefficient R between the mean relative prediction error and the number of parallel TCP flows is less than 0.1.

C. Outcome

Our experimental results have shown how, using the model derived in this section, one can effectively predict the throughput of parallel TCP for a wide range of parallelism relying only

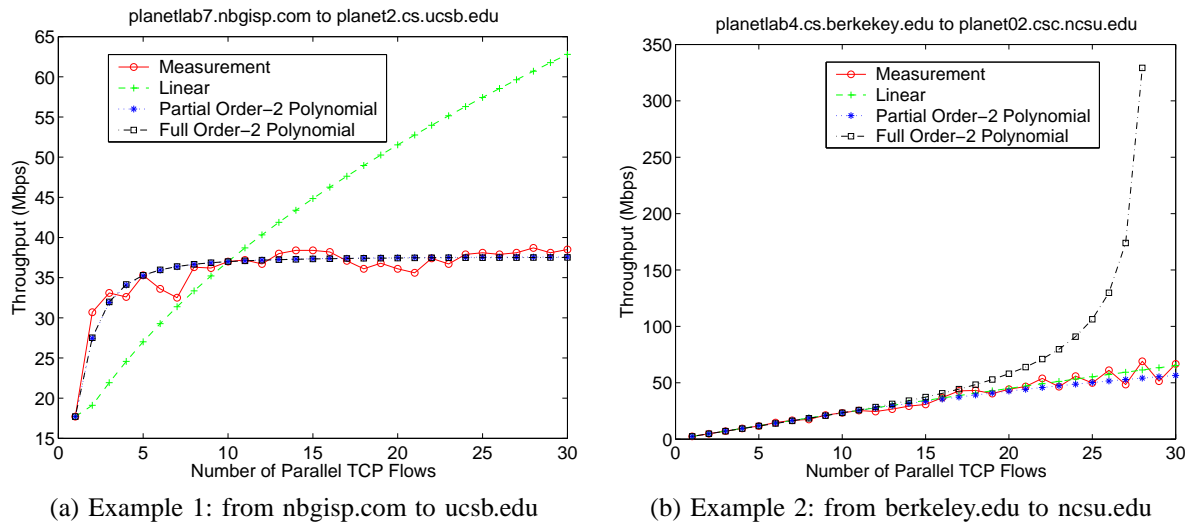


Fig. 9. Parallel TCP throughput prediction.

Source Host	Destination Host	Mean	Standard Deviation
planetlab2.postel.org	planetlab-1.cmcl.cs.cmu.edu	0.0822	0.1264
planetlab2.it.uts.edu.au	planetlab1.diku.dk	0.1110	0.4061
planetlab1.millennium.berkeley.edu	planetlab3.cs.uoregon.edu	0.0910	0.1325
planetlab-1.it.uu.se	planet2.cs.ucsb.edu	0.1302	0.2599
ds-pl1.technion.ac.il	planetslug3.cse.ucsc.edu	-0.0453	0.1084
planetlab-2.cs.princeton.edu	planetlab9.cs.berkeley.edu	0.07049	0.15274
planetlab2.flux.utah.edu	planet1.cs.ucsb.edu	-0.1081	0.2583
planetlab2.chin.internet2.planet-lab.org	planetlab5.millennium.berkeley.edu	0.0211	0.0594
planetlab2.frankfurt.interxion.planet-lab.org	planet1.cc.gt.atl.ga.us	-0.1676	0.2887
planetlab2.bgu.ac.il	planetlab1.it.uts.edu.au	-0.0533	0.1264
planetlab2.cs.berkeley.edu	planetslug2.cse.ucsc.edu	-0.0138	0.1033
planet1.calgary.canet4.nodes.planet-lab.org	planetlab2.sanjose.equinix.planet-lab.org	-0.1088	0.1331
planetlab2.cs-ipv6.lancs.ac.uk	planetlab4.cs.berkeley.edu	-0.0518	0.1514
planetlab2.cs.uoregon.edu	planet02.csc.ncsu.edu	-0.0522	0.0895
planetlab2.postel.org	planetlab2.it.uts.edu.au	-0.1208	0.2825
planetlab2.flux.utah.edu	planetlab2.chin.internet2.planet-lab.org	0.1235	0.4334
planetlab2.frankfurt.interxion.planet-lab.org	planetlab2.bgu.ac.il	-0.0956	0.2232
planetlab9.millennium.berkeley.edu	planetlab2.cs.berkeley.edu	-0.0070	0.0111
planetlab2.cs-ipv6.lancs.ac.uk	planetlab2.cs.uoregon.edu	-0.2581	0.1828
planetlab-2.it.uu.se	planetlab3.sanjose.equinix.planet-lab.org	0.0664	0.0848
pli1-pa-3.hpl.hp.com	planetlab7.millennium.berkeley.edu	0.0400	0.1542
planetlab1.cs.berkeley.edu	planetlab01.ethz.ch	-0.1546	0.1937
planetlab3.cs.uoregon.edu	planetlab02.cs.washington.edu	-0.1346	0.1068
planetlab7.nbgisp.com	planet2.cs.ucsb.edu	0.0033	0.0394
planetslug3.cse.ucsc.edu	planetlab9.cs.berkeley.edu	-0.2750	0.2743
planet1.cs.ucsb.edu	planetlab5.millennium.berkeley.edu	0.0743	0.2118
planet1.cc.gt.atl.ga.us	planetlab1.it.uts.edu.au	-0.1753	0.3964
planetlab4.millennium.berkeley.edu	planetslug2.cse.ucsc.edu	5.1483e-04	0.0028
planetlab4.cs.berkeley.edu	planet02.csc.ncsu.edu	-0.0136	0.0882
planetlab2.postel.org	planetlab2.cs.berkeley.edu	0.0029	0.1051
planetlab-2.cmcl.cs.cmu.edu	planetlab2.cs.uoregon.edu	0.0258	0.0664
planetlab2.flux.utah.edu	planetlab3.sanjose.equinix.planet-lab.org	0.1742	0.1491
planetlab2.frankfurt.interxion.planet-lab.org	planetlab2.tau.ac.il	0.03886	0.2650
planetlab9.millennium.berkeley.edu	planetlab1.flux.utah.edu	0.0922	0.1430
planetlab2.cs-ipv6.lancs.ac.uk	planetlab7.millennium.berkeley.edu	-0.1643	0.1345
planetlab-2.it.uu.se	planetlab1.enel.ucalgary.ca	-0.1604	0.1833
s2_803.ie.cuhk.edu.hk	planetlab01.ethz.ch	-0.0375	0.5193
planet2.pittsburgh.intel-research.net	planetlab02.cs.washington.edu	0.190	0.4300
planetlab2.millennium.berkeley.edu	planetlab1.it.uts.edu.au	-0.1769	0.1695
planetlab1.cs.berkeley.edu	planetslug2.cse.ucsc.edu	0.0200	0.0912
planetlab7.nbgisp.com	planetlab5.millennium.berkeley.edu	0.0093	0.0747

Fig. 10. Relative Prediction Error Statistics for Parallel TCP Throughput.

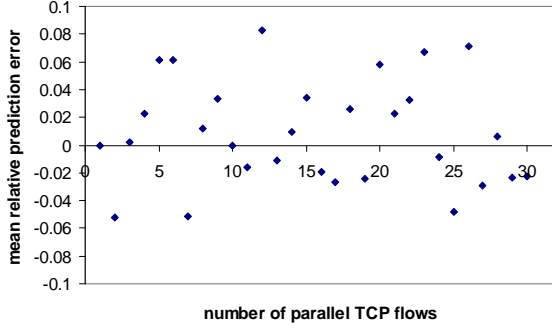


Fig. 12. Relative prediction error for parallel TCP throughput as a function of number of parallel TCP flows.

on two active probes at different levels of parallelism. In the following section we try to estimate the effect of parallel TCP in the existing cross traffic for a given level of parallelism, the last “piece” necessary to make the `TameParallelTCP()` call possible.

V. TAMING PARALLEL TCP

There are a number of considerable challenges when trying to estimate the effect on cross traffic with an online system running on the end points:

- The available bandwidth on the bottleneck link(s) is unknown.
- The cross traffic on the bottleneck link(s) (the offered load) is unknown.
- Making use of an additional network measurement tool (such as Pathload [19], [20]) to determine the current load on the path is problematic since it can take a long time to converge. In addition, the measurement accuracy cannot be guaranteed. One would like to avoid any additional overhead beyond the required two active probes necessary to predict the throughput of parallel TCP flows.

In what follows, we make simplifying assumptions about the cross traffic’s view of the shared links on the path in order to provide an estimate of impact on the cross traffic from the same two probes from which we derived the throughput curve in the previous section.

A. Algorithm

We assume that all TCP connections, including our parallel TCP flows and the cross traffic, share the same loss rate on a bottleneck link. This assumption is valid if at least one of the two following conditions is satisfied:

- The cross traffic has an RTT similar to our parallel TCP flows. In that case, all connections are very likely to have their congestion window synchronized, and thus share the same loss rate. This fact has been independently verified by other research groups [33], [30], [31].
- The router on the bottleneck link is using Random Early Detection (RED) [14] as its queue management policy, something that is becoming increasingly more common. Research has demonstrated that with RED, different flows

roughly experience the same loss rate (the RED rate, which depends on the queue occupancy) under steady state [14], [31].

Our approach to determining the effect of parallel TCP on cross traffic is based on our algorithm to estimate the parallel TCP throughput (Section IV). The key idea is to estimate $p_n \times RTT_n^2$ as a function of the number of parallel TCP flows. Based on the assumption that cross traffic shares the same loss rate as parallel TCP flows, we can then use the simple TCP throughput model (Equation 1) to estimate the relative change to the cross traffic throughput.

Recall in section IV we model $p_n \times RTT_n^2$ with a partial order 2 polynomial function $a \times n^2 + b$ (Equation 5). After having obtained the two necessary measurements, we can calculate the value of a and b and are now able to estimate the loss rate as a function of the number of parallel TCP flows.

Relying on our assumptions, we have also obtained the loss rate of the cross traffic as a function of the number of parallel TCP flows n given there are m cross traffic flows (recall that we also assume that m is relatively stable).

Thus, based on Equation 1, we can now estimate the relative change on each of the individual TCP throughputs without knowing m using the following equation:

$$relc = \frac{\frac{MSS \times C}{RTT_{n1} \times \sqrt{p_{n1}}} - \frac{MSS \times C}{RTT_{n2} \times \sqrt{p_{n2}}}}{\frac{MSS \times C}{RTT_{n1} \times \sqrt{p_{n1}}}} \quad (15)$$

$$= 1 - \frac{RTT_{n1} \times \sqrt{p_{n1}}}{RTT_{n2} \times \sqrt{p_{n2}}} \quad (16)$$

Here, $relc$ is the relative throughput change for each flow, MSS and C are constants as described in Section IV, and $RTT_n \times \sqrt{p_n}$ can be estimated according to the number of parallel TCP flows n as described above.

B. Evaluation

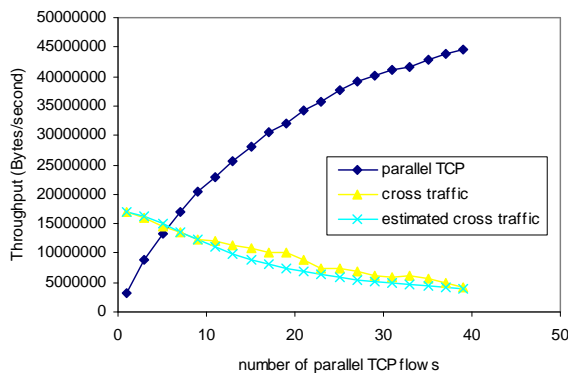
We have done a thorough ns2-based evaluation of our cross traffic estimator. Simulation experiments allow us to analyze our estimator under controlled, reproducible settings including bottleneck bandwidth and cross traffic characteristics.

Our simulation configuration was already introduced in Section III. We consider the same set the scenarios presented there. As in Section III, we employ Qiu et al’s [30] simulation topology (Figure 1). The representativeness of this topology was discussed in previous research [30], [31].

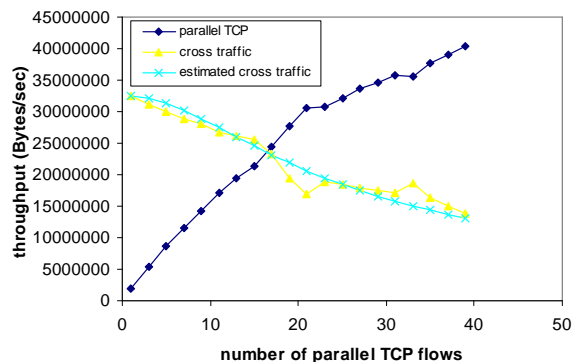
Figure 13 shows two examples, for Scenarios 4 and 6, of the performance of our estimator. In these cases we can accurately predict the impact on cross traffic as a function of the parallelism level using only two probes, the same probes we use to predict the throughput of the parallel flows as a function of parallelism level.

We summarize our prediction results as a CDF of the relative error in predicting the impact on cross traffic across all of our scenarios in Figure 14. We can see that 90% of predictions have relative prediction error less than 0.25. The cross traffic estimator is slightly biased. It conservatively predicts a greater impact on the cross traffic on average.

To further evaluate our cross traffic estimation algorithm, we designed a more complicated topology with two groups of



(a) Scenario 4 with 5 cross traffic



(b) Scenario 6 with 15 cross traffic

Fig. 13. Examples of cross traffic estimation with simulations in section IV.

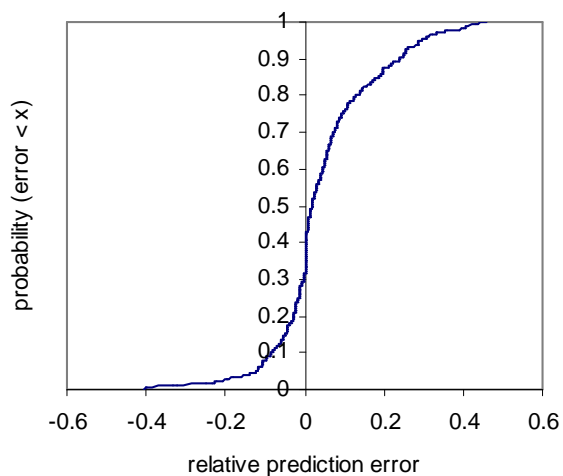


Fig. 14. Cumulative distribution function of relative prediction error for cross traffic estimation for all the simulations with 6 scenarios as described in Figure 2.

cross traffic. The topology and the simulation configuration is shown in Figure 15. Each simulation is 100 seconds long with cross traffic starting randomly between 0 and 8 seconds and all the parallel TCP flows starting at 10 seconds. We applied our same algorithm for estimation of cross traffic. The results are presented in Figure 16, which clearly shows the effectiveness of our approach.

We also tested the cross traffic estimator for scenarios in which different TCP flows have different RTTs, and where RED is not used on the routers. Our estimator shows the right trend of the cross traffic throughput change, though accurate prediction cannot be guaranteed as flows with longer RTT tend to have higher loss rate than parallel TCP flows and vice versa. In essence, in situations in which cross traffic RTT and loss rate is unknown, our estimator is less accurate.

C. Outcome

In this section, we have demonstrated the feasibility of predicting the impact on cross traffic of a parallel TCP transfer as a function of the degree of parallelism. Under the assumption that all flows share the same loss rate, we

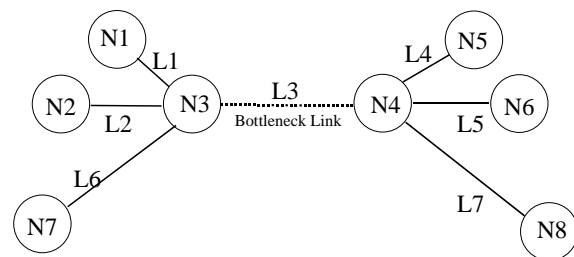


Fig. 15. Simulation topology for further evaluations of cross traffic estimation. L1 and L4 have latency 3ms, L2 and L5 have latency 6ms, L6 and L7 have latency 10ms. L3 have latency 50ms and bottleneck bandwidth 1000Mbps. N3 is using RED queue management policy. Parallel TCP flows go from N2 to N6. Cross traffic group 1 goes from N7 to N8. Cross traffic group 2 goes from N1 to N5.

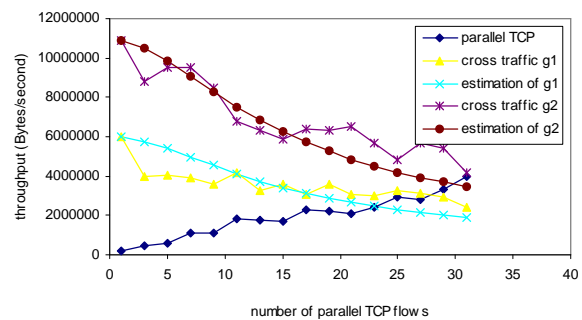


Fig. 16. Estimation results with 14 TCP flows in cross traffic group 1 (g1) and 14 TCP flows in cross traffic group 2 (g2).

can accurately predict the relative impact using the same two measurement probes used to predict the throughput of the parallel TCP transfer as a function of the degree of parallelism.

Combining these two predictions, we can implement the `TameparallelTCP()` API call:

- 1) Execute two probes at different parallelism levels. (1, 5) appears to work fine.
- 2) Using the probe results, estimate the parallel TCP throughput as a function of the number of parallel TCP flows n using the techniques of the previous section.
- 3) Using the probe results, estimate the relative impact on cross traffic as a function of n using the techniques of this section. and the relative throughput

- 4) Conduct a binary search on the cross traffic impact function, looking for the degree of parallelism, l , that has the largest impact less than that permitted in the API call.
- 5) Return l , and the impact and throughput predictions at parallelism l .

The cost of this implementation is dominated by executing the two probes.

VI. CONCLUSIONS AND FUTURE WORK

We have shown how to predict both parallel TCP throughput and its impact on cross traffic as a function of the degree of parallelism using only two probes at different parallelism levels. Both predictions are monotonically increasing with parallelism level. Hence, the `TameParallelTCP()` function can be implemented using a simple binary search.

We have made a number of simplifying assumptions about the cross traffic in order to predict impact on it while having no knowledge of the actual cross traffic. While these assumptions are reasonable in many cases, we are now working on how to relax them. We will soon make available a straightforward library implementation of our `TameParallelTCP()` API call.

REFERENCES

- [1] <http://dast.nlanr.net/projects/iperf/>.
- [2] <http://www.isi.edu/nsnam/ns/>.
- [3] <http://www.planet-lab.org>.
- [4] <http://www.web100.org>.
- [5] ALLCOCK, W., BESTER, J., BRESNAHAN, J., CERVENAK, A., LIMING, L., AND TUECKE, S. GridFTP: Protocol extensions to ftp for the grid. Tech. rep., Argonne National Laboratory, August 2001.
- [6] ALLCOCK, W., BESTER, J., BRESNAHAN, J., CHERVENAK, A., FOSTER, I., KESSELMAN, C., MEDER, S., NEFEDOVA, V., QUESNEL, D., AND TUECKE, S. Data management and transfer in highperformance computational grid environments. *Parallel Computing* 28 (2002).
- [7] BOLLIGER, J., GROSS, T., AND HENGARTNER, U. Bandwidth modeling for network-aware applications. In *INFOCOM (3)* (1999), pp. 1300–1309.
- [8] CARTER, R., AND CROVELLA, M. Measuring bottleneck link speed in packet-switched networks. *Performance Evaluation*, 28 (1996), 297–318.
- [9] DOVROLIS, C., RAMANATHAN, P., AND MOORE, D. What do packet dispersion techniques measure? In *INFOCOM* (2001), pp. 905–914.
- [10] DOWNEY, A. B. Using pathchar to estimate internet link characteristics. In *Measurement and Modeling of Computer Systems* (1999), pp. 222–223.
- [11] EGGERT, L., HEIDEMANN, J., AND TOUCH, J. Effects of ensemble-tcp.
- [12] FALL, K., AND FLOYD, S. Simulation-based comparisons of Tahoe, Reno and SACK TCP. *Computer Communication Review* 26, 3 (July 1996), 5–21.
- [13] FISK, M., AND FENG, W. Dynamic right-sizing: Tcp fw-control adaptation. In *Supercomputing (SCO1)* (2001).
- [14] FLOYD, S., AND JACOBSON, V. Random early detection gateways for congestion avoidance. *IEEE/ACM Transactions on Networking* 1, 4 (1993), 397–413.
- [15] FOSTER, I. Globus web page. Tech. Rep. <http://www.mcs.anl.gov/globus>, Argonne National Laboratory.
- [16] HACKER, T., ATHEY, B., AND NOBLE, B. The end-to-end performance effects of parallel tcp sockets on a lossy wide-area network. In *16th IEEE/ACM International Parallel and Distributed Processing Symposium (IPDPS)* (2002).
- [17] HACKER, T. J., NOBLE, B. D., AND D., B. The effects of systemic packet loss on aggregate tcp fwbs. In *IEEE/ACM Supercomputing* (2002).
- [18] HU, N., AND STEENKISTE, P. Evaluation and characterization of available bandwidth probing techniques. *IEEE JSAC Special Issue in Internet and WWW Measurement, Mapping, and Modeling* 21, 6 (August 2003).
- [19] JAIN, M., AND DOVROLIS, C. End-to-end available bandwidth: Measurement methodology, dynamics, and relation with tcp throughput. In *ACM SIGCOMM* (2002).
- [20] JAIN, M., AND DOVROLIS, C. Pathload: A measurement tool for end-to-end available bandwidth. In *Passive and Active Measurement Workshop* (2002).
- [21] JIN, G., YANG, G., CROWLEY, B., AND AGARWAL, D. Network characterization service (ncs). In *10th IEEE Symposium on High Performance Distributed Computing, Aug. 2001*. (2001).
- [22] LAI, K., AND BAKER, M. Nettimer: A tool for measuring bottleneck link bandwidth. In *USENIX Symposium on Internet Technologies and Systems* (2001), pp. 123–134.
- [23] LEE, J., GUNTER, D., TIERNEY, B., ALLCOCK, B., BESTER, J., BRESNAHAN, J., AND TUECKE, S. Applied techniques for high bandwidth data transfers across wide area networks. In *International Conference on Computing in High Energy and Nuclear Physics, Beijing, China, September 2001*.
- [24] LOWEKAMP, B., MILLER, N., SUTHERLAND, D., GROSS, T., STEENKISTE, P., AND SUBHLOK, J. A resource monitoring system for network-aware applications. In *Proceedings of the 7th IEEE International Symposium on High Performance Distributed Computing (HPDC)* (July 1998), IEEE, pp. 189–196.
- [25] LU, D., QIAO, Y., DINDA, P. A., AND BUSTAMANTE, F. E. Characterizing and predicting tcp throughput on the wide area network. Tech. Rep. NWU-CS-04-34, Northwestern University, Department of Computer Science, 4, 2004.
- [26] MATHIS, M., HEFFNER, J., AND REDDY, R. Web100: Extended tcp instrumentation for research, education and diagnosis. *ACM Computer Communications Review* 33, 3 (July 2003).
- [27] MATHIS, M., SEMKE, J., AND MAHDAVI, J. The macroscopic behavior of the tcp congestionavoidance algorithm. *Computer Communication Review* 27, 3 (1997).
- [28] MOGUL, J., AND DEERING, S. A framework for defining empirical bulk transfer capacity metrics, rfc3148, November 1990.
- [29] PADHYE, J., FIROIU, V., TOWSLEY, D., AND KUROSE, J. Modeling tcp throughput: A simple model and its empirical validation. In *ACM SIGCOMM* (1998).
- [30] QIU, L., ZHANG, Y., AND KESHAV, S. On individual and aggregate TCP performance. In *ICNP* (1999), pp. 203–212.
- [31] QIU, L., ZHANG, Y., AND KESHAV, S. Understanding the performance of many TCP fwbs. *Computer Networks (Amsterdam, Netherlands: 1999)* 37, 3–4 (2001), 277–306.
- [32] RIBEIRO, V., RIEDI, R., BARANIUK, R., NAVRATIL, J., AND COTRELL, L. pathchirp: Efficient available bandwidth estimation for network paths. In *Passive and Active Measurement Workshop* (2003).
- [33] SHENKER, S., ZHANG, L., AND CLARK, D. Some observations on the dynamics of a congestion control algorithm. *ACM Computer Communication Review* (1990).
- [34] SIVAKUMAR, H., BAILEY, S., AND GROSSMAN, R. L. Pockets: The case for application-level network striping for data intensive applications using high speed wide area networks. In *Supercomputing* (2000).
- [35] SMITH, F. D., HERNANDEZ-CAMPOS, F., JEFFAY, K., AND OTT, D. What TCP/IP protocol headers can tell us about the web. In *SIGMETRICS/Performance* (2001), pp. 245–256.
- [36] STRAUSS, J., KATABI, D., AND KAASHOEK, F. A measurement study of available bandwidth estimation tools. In *Internet Measurement Conference* (2003).
- [37] TIERNEY, B. Tcp tuning guide for distributed application on wide area networks. *USENIX & SAGE Login* 26, 1 (2001).
- [38] ZHANG, Y., BRESLAU, L., PAXSON, V., AND SHENKER, S. On the Characteristics and Origins of Internet fw rates. In *ACM SIGCOMM* (2002).
- [39] ZHANG, Y., DUFFIELD, N., PAXSON, V., AND SHENKER, S. On the constancy of internet path properties. In *ACM SIGCOMM Internet Measurement Workshop* (November 2001).