

# A Multipath Background Network Architecture

Ravi Kokku<sup>†</sup> Aniruddha Bohra<sup>†</sup> Samrat Ganguly<sup>†</sup> Arun Venkataramani<sup>‡</sup>

<sup>†</sup>NEC Laboratories America Inc., Princeton, NJ    <sup>‡</sup>University of Massachusetts, Amherst, MA  
{ravik, bohra, samrat}@nec-labs.com, arun@cs.umass.edu

*Abstract*—Background transfers, or transfers that humans do not actively wait on, dominate the Internet today. In today’s best-effort Internet, background transfers can interfere with foreground transfers causing long wait times, thereby hurting human productivity. In this paper, we present the design and implementation of a background network, **Harp**, that addresses this problem. **Harp** has three significant advantages over recent end-host based background transport protocols; **Harp** (i) uses multiple paths to exploit path diversity and load imbalance in the Internet to tailor network resource allocation to human needs, (ii) provides better fairness and utilization compared to unipath end-host protocols, and (iii) can be deployed at either end-hosts or enterprise gateways, thereby aligning the incentive for deployment with the goals of network customers. Our evaluation using simulations and a prototype on Planetlab suggests that **Harp** improves foreground TCP transfer time by a factor of five and background transfer time by a factor of two using just two extra paths per connection.

## I. INTRODUCTION

Background transfers, or transfers that humans do not actively wait on, dominate Internet traffic today. Examples include peer-to-peer file swarming [1, 2], content distribution [3, 4], Web [5, 6] and media [7] prefetching, remote backup [8, 9], software updates [10] etc. A recent study [11] suggests that up to 90% of Internet traffic may be background. However, today’s best-effort Internet does not distinguish between background and foreground (or interactive) transfers. As a result, background transfers *interfere* with foreground transfers and inconvenience humans, e.g., online purchases or important email can get delayed when an enterprise network is bogged down by large software updates; or the quality of video-conferencing can degrade because of concurrent file swarming or file system backup traffic.

Two approaches exist to address interference: (i) priority queuing at routers [12], and (ii) end-host based transport protocols [13, 14, 15] that provide lower-than-best-effort service without support from routers. Both approaches suffer from two fundamental problems.

- *Arbitrary delay to background transfers*: Prioritization using either approach can arbitrarily delay background transfers. Furthermore, end-host based protocols are imperfect due to delayed feedback and suffer from a marked drop in utilization to achieve near-zero interference. So, users have to plan their activities around longer download times for such transfers. For example, even though users do not actively wait on a movie download, a download that takes two hours instead of one can make users unhappy.
- *Deployability*: Priority queuing at routers appears impractical due to the lack of architectural consensus today [16]. On the other hand, modifying the applications and/or op-

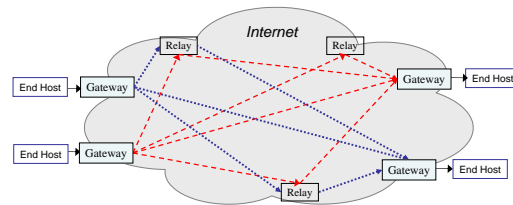


Fig. 1. Harp’s basic architecture with gateways and in-network relays.

erating systems at end-hosts (although easier than modifying routers) is still a significant barrier to widespread adoption, and is worsened by the reluctance of users to install software that makes them “good network citizens” by slowing down their background transfers.

In this paper, we present the design and implementation of a network architecture, **Harp**, for a background-dominated Internet, which addresses the above problems using two ideas.

- *Wait-time as a first-class metric*: **Harp** prioritizes foreground over background transfers to reduce active wait-time. More importantly, **Harp** uses multipath routing to dissipate background transfers to less-utilized regions of the network. Thus, **Harp** exploits the path diversity, load imbalance, and spare capacity in the Internet to ensure quicker completion of both foreground and background transfers.

Our algorithmic contribution is a background multipath congestion controller that jointly optimizes rate control as well as routing to improve utilization and fairness of network-wide spare resources. More importantly, our effort to build a practical system revealed a shortcoming of previously proposed multipath controllers for foreground transport [17, 18], namely, the utilization is significantly below that predicted theoretically when the number of flows is small. Our algorithm addresses this problem by smoothly transitioning between joint and independent control of rate and routing depending on the number of other extant flows on the path. **Harp**’s controller is applicable to foreground transport as well with similar utilization benefits.

- *Deployment as an edge service*: **Harp** supports background transfers as an edge-service. **Harp** is implemented as an overlay network consisting of two types of transit nodes: *relays* and *gateways*. In the envisioned deployment scenario (Figure 1), an edge-service provider supports relays distributed across the Internet to enable path diversity, and organizations install gateways at their network exits—a popular edge-service business model [19, 20] to-

day. Harp can also be deployed directly at end-hosts and use peer nodes as relays, obviating transit nodes. However, our position is that a practical architecture must align the incentive for deployment with the goals of network customers that are often physically disjoint from end-users as in the case of enterprises, public access points, overlay network services etc. Thus, a gateway is a sweet spot to incent deployment compared to modifying routers or diverse end-hosts of unwilling users.

Our systems contribution is a prototype of Harp built using Click [21] and deployed on Planetlab [22]. We address non-trivial engineering challenges to enable background multipath congestion control as an edge-service without support from end-hosts or routers. Our extensive evaluation using this prototype as well as *ns* simulations demonstrate that Harp can significantly improve transfer completion times. For example, our simulations show that Harp improves foreground TCP transfer completion time by a factor of five. Our prototype experiments show that background completion time can be reduced by a factor of two using just two extra paths per connection.

The rest of the paper is organized as follows. Sections II and III present the design and implementation of Harp. Section IV presents the results of *ns* and prototype-based experiments. Section V presents related work and Section VI concludes.

## II. SOLUTION OVERVIEW

Harp uses a multipath network for background transfers deployed as an edge-service and leaves foreground transfers unchanged. Foreground transfers are those that humans actively wait on, whereas during background transfers, humans typically engage themselves in other activities. We note that categorizing all transfers into two classes is a simplifying assumption, but is appropriate for a background-dominated Internet. Harp does not employ multipath routing for foreground transfers as they are typically too short-lived to get steady-state utilization or fairness benefits.

The design and implementation of Harp is guided by the following principles.

- 1) *Wait-time first*: Noninterference of background transfers must be treated as a safety property to reduce active wait-time. For example, a loss of a foreground SYN or data packet in slow start can severely hurt its completion time. Nevertheless, a short foreground transfer should minimally disrupt a background flow in steady-state. Harp’s congestion controller accomplishes both objectives.
- 2) *Evolvability*: An evolvable solution must not embed mechanisms into routers to optimize today’s background-dominated Internet. Harp is implemented as an edge-service without support from routers.
- 3) *Incentive compatibility*: Unlike broadband users who might use Harp to reduce wait-time across their own transfers, users at public access points, enterprises, and overlay network services do not have an incentive to slow down their background transfers. Harp’s gateway obviates any modification to or cooperation from end-hosts.
- 4) *Fairness*: An architecture must preserve fairness of resource allocation. Harp’s multipath congestion controller

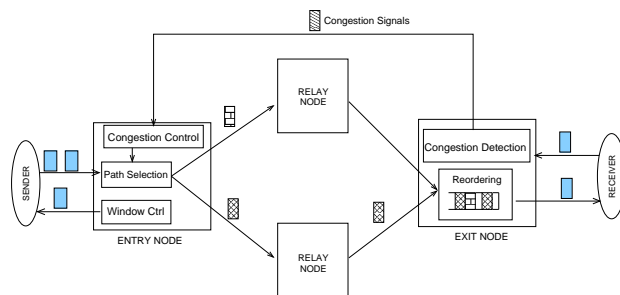


Fig. 2. Harp Architecture

controls the throughput achieved by a flow on each path so as to achieve TCP-fairness across all flows using Harp.

- 5) *End-to-end semantics*: An edge-service should be seamless to legacy applications and operating systems that are in charge of end-to-end reliability. Harp’s transit nodes unreliaibly buffer, congestion control, and forward data packets. Acknowledgments and retransmissions are handled by the end-host, preserving end-to-end reliability semantics.

Figure 2 illustrates the basic sequence of steps for packet delivery in Harp. Packets from the end-host enter the entry gateway and go to the exit gateway either directly or through one of the relay nodes. The exit gateway reorders packets and delivers them to the receiving end-host. The acknowledgments generated by the receiver for all packets are sent back directly from the exit to the entry gateway. The multipath background congestion controller resides at the entry and exit gateways.

## III. DESIGN

Harp uses a multipath controller to dissipate background traffic to less-utilized parts of the network, out of the way of foreground transfers. Harp’s controller builds upon Kelly and Voice’s (KV) multipath routing and congestion controller [17] and TCP Nice [13], a unipath background congestion controller. Similar to KV, Harp shares congestion information across multiple paths in a connection to achieve a fair allocation of network-wide resources. Similar to Nice, Harp uses delay-based signals to proactively detect congestion, and aggressively backs off in response. The rest of this section presents the congestion control algorithm, and the challenges and design alternatives in realizing it without end-host support.

### A. Multipath Congestion Control

The congestion control algorithm executes on the entry gateway. Given a set of paths on which packets can be sent, the algorithm estimates the size of the congestion window on each path that indicates a stable sending rate on the path. Our algorithm uses the following multiplicative increase multiplicative decrease (MIMD) scheme —(1) MI: on each positive acknowledgment on a path, increment the congestion window of the path by an MI parameter  $\alpha$ , and (2) MD: on a congestion signal on a path, decrement the window of the path by an MD parameter  $\beta$  times the weighted sum of congestion window on the current path ( $w_i$ ) and the total congestion window on all paths ( $W$ ).

Algorithm 1 shows the pseudo code of the above congestion control behavior. Line 6 represents the MI part and Lines 9

---

**Algorithm 1** Harp's Congestion Control Algorithm

---

```
1: Definitions:  $i$ : path (numbered from 1 to  $n$ )
2:  $w_i$ : congestion window (in bytes) on path  $i$ 
3:  $W$ : total congestion window =  $\sum_1^n w_i$ 
4:
5: On ack for path  $i$ 
6:  $w_i \leftarrow w_i + \alpha$ 
7:
8: On delay signal for path  $i$ 
9:  $w_i \leftarrow \max(1, w_i - \beta \times (w_i \times \xi + W \times (1 - \xi)))$ 
10:
11: On loss signal for path  $i$ 
12:  $w_i \leftarrow \max(1, w_i - W/2)$ 
```

---

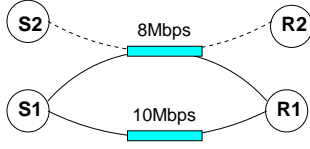


Fig. 3. Example of two senders sharing one path

and 12 represent the MD part. We use both increased packet delays and packet losses as signals of congestion. To prevent the algorithm from over-reacting to congestion indications and thus ensure stability, we let the congestion window be decremented only once per round-trip time on a path after observing a threshold number of congestion indications. Further, the window is not incremented within one round-trip time of the previous decrement.

Setting the value of  $\xi$  to different values allows us to instantiate several variants of multipath congestion control.

- **Independent:**  $\xi = 1$  makes the multiplicative decrease on a path proportional to the sending rate on that path. This is same as independent congestion control on each path, i.e., each subflow on a path operates as an individual TCP flow [23]. Hence, a multipath connection using such independent control with  $n$  paths operates as  $n$  TCP flows.
- **Joint:**  $\xi = 0$  makes the multiplicative decrease part same as that of joint routing and congestion control proposed by Han et al. [18] and Kelly et al. [17]. With such control, each multipath connection behaves as a single TCP flow on its best set of paths (one or more).

Following either of the above two approaches has both advantages and disadvantages. When multiple background transfers share a set of paths, Independent can be unfair to transfers that have fewer paths at their disposal. For instance, consider Figure 3. The sender-receiver pair S1R1 has access to two paths, whereas S2R2 only has one path. When using independent, S1R1 gets an equal share of bandwidth on the path shared with S2R2 that leads to significantly higher overall bandwidth to S1R1 (14Mbps) than S2R2 (4Mbps). Using Joint instead, would move each transfer to its best set of paths, i.e. S1R1 occupies only the 10Mbps path and S2R2 completely occupies 8Mbps path, thereby leading to a fairer allocation of resources.

Using Joint when S1R1 is the only transfer, however, can lead to under-utilization of a path. This is because even if S1R1 is the only transfer on a path, the transfer backs-off proportional

to the *cumulative* congestion window on all paths thereby losing more throughput than Independent control. This behavior of Joint could be unnecessarily aggressive for background transfers whose goal is to reap as much spare capacity as possible.

To get the best of both Independent and Joint, we explore a variant that combines the behaviors by adapting the value of  $\xi$ .

- **Adaptive:** We set  $\xi = \frac{w_i}{M_i}$ , where  $w_i$  is the congestion window during decrement, and  $M_i$  is the maximum congestion window size observed by the transfer on path  $i$ .

Adaptive control has the following properties. When a multipath connection is the only one active on one or more of its paths, the multiplicative decrement on such paths behaves more like Independent because  $w_i$  is close to  $M_i$ . As the number of transfers sharing the path increases, the behavior of Adaptive becomes closer to Joint. However, for this technique to work well, each transfer should observe the maximum congestion window  $M_i$  at some time during its activity. Due to this requirement, Adaptive could take longer to achieve the same fairness property as Joint. We believe that this tradeoff is acceptable for long running transfers.

We formally summarize the above observations in the theorem below (and present a proof in the technical report [24]).

*Theorem 1: (a) Joint control of congestion and routing can degrade steady-state throughput by an  $\Omega(\ln(C))$  factor compared to independent control along a path of bottleneck capacity  $C$ , (b) Adaptive control achieves the same steady-state throughput as independent when there is exactly one flow along the path.*

### B. Path Selection

On the arrival of a packet from the sender, the entry gateway attempts to balance load across the available paths by choosing for each packet a path with minimum  $\frac{\text{bytes\_in\_nw}_i}{\text{cwnd}_i}$ , where  $\text{bytes\_in\_nw}_i$  represents the number of unacknowledged bytes sent on path  $i$ . The same expression is used to stripe packets even when  $\text{bytes\_in\_nw}_i$  exceeds  $\text{cwnd}_i$  to ensure load balancing on the paths. Once the entry gateway selects a path (i.e., the corresponding relay node) to send the packet through, it uses packet encapsulation to facilitate routing packets through the relay node. Additionally, each encapsulated packet carries a multipath header that contains a packet type (to represent data, probe, loss signal, etc.), a timestamp (representing when the packet left the entry gateway) and a path identifier (which is the IP address of the relay node). This header allows the exit gateway to identify the path taken by the packet, and detect and associate congestion to the path based on the delay observed by the packet.

### C. Reordering

Since packets are sent on several paths with different latencies, they can arrive out of order and cause the receiver to send duplicate acknowledgments to the sender. Such duplicate acknowledgments falsely indicate packet loss to the sender that can lead to a substantial reduction of its congestion window thereby reducing its throughput. To avoid such false indications, our mechanism reorders packets received from multiple paths at the exit gateway before sending them to the receiver

using the sequence numbers in packets’ TCP headers. If packets are in sequence, they are sent immediately to the receiver. Otherwise, they are kept in a reorder queue till either the packets before them arrive (and thus complete the sequence) or a timer expires. Since the reorder delay required at the exit gateway is governed by the path with the longest delay, we set the timer to a value that is a factor  $\rho$  of the minimum delay on the longest path.

#### D. Congestion Indication

To reduce interference, background transfers should *back-off* as early as possible when required. We achieve this behavior by reacting to increased packet delays, in addition to losses.

*Detecting Increased Delays:* Harp uses relative increase in one-way packet delays as early indicators of congestion [14]. To achieve this, each packet is timestamped (in the multipath header) at the entry gateway before being sent on a path. The exit gateway calculates the one-way delay using the timestamp and its current time. Also, the exit gateway keeps track of minimum ( $dmin_i$ ) and maximum ( $dmax_i$ ) delays observed by packets on each path for a connection. If a packet observes a delay greater than  $\Delta_i = dmin_i + (dmax_i - dmin_i) \times \delta$ , where  $\delta$  is a threshold parameter, the exit gateway sends a congestion indication to the entry gateway. Observe that in the above condition, we only use the relative difference in packet delays to infer congestion; hence the clocks on the entry and exit gateways need not be perfectly synchronized. Congestion can be indicated by either sending an explicit message to the entry gateway or by piggybacking the indication on a returning acknowledgment. We choose the latter in Harp for efficiency.

*Detecting Losses:* Since the exit gateway maintains a reorder queue, it can detect packet losses earlier than the receiver. The exit gateway maintains a variable  $last\_byte\_rcvd_i$  for each path  $i$  that indicates the highest byte received on the path. Additionally, it maintains a  $rcvnext$  variable per connection that indicates the next byte expected *in sequence*. When the  $last\_byte\_rcvd_i$  on each path exceeds  $rcvnext$ , the exit gateway detects a possible packet loss, and sends a loss indication to the entry gateway. The message also contains the range of missing bytes, which can be determined by  $rcvnext$  and the sequence number of the packet at the head of the reorder queue. Observe that the exit gateway can not exactly determine which path the loss occurred. Hence the range is used at the entry gateway to determine the path(s) on which the packets containing the missing bytes were sent. For each path on which any of the missing bytes were sent, the congestion window is reduced as described in Section III-A. This technique of detecting packet losses is simpler and quicker than waiting for and interpreting duplicate acknowledgments from the receiver.

As an optimization, we could cache unacknowledged packets at the entry gateway and respond to loss indications by retransmitting packets and thus avoid exposing the loss to the sender. To avoid exposing the loss, we would also need to interpret and suppress duplicate acknowledgments from the receiver. Note that doing so would not violate our first guideline of maintaining end-to-end semantics described in Section II. However, we

do not currently incorporate this optimization in our prototype.

Formally, we state the following theorem on interference caused by Harp, and present a proof in the technical report [24].

*Theorem 2: The interference, or reduction in throughput of foreground flows, that Harp’s multipath background controller inflicts is bounded by a factor  $O(e^{(-\frac{B}{m} \cdot (1-\delta)(1-\beta)})}$ , where  $B$  is the bottleneck buffer capacity,  $\delta$  the delay threshold,  $\beta$  the MD parameter, and  $m$  the number of foreground flows along the path, independent of the number of background flows.*

#### E. Congested-path Suppression

Congested paths can increase the packet delays substantially and even cause losses thereby causing packets on other paths to wait in the reorder queue. Such a wait degrades the throughput of the connection because new packets are not sent by the sender unless the receiver acknowledges the old packets.

To reduce the impact of such congested paths on throughput, whenever the congestion window for a path reaches MINCWND, the path is temporarily marked as choked. No subsequent packets are sent on this path unless it is unchoked. From then on, a probe (unchoke request) with timestamp is sent on each choked path periodically. If the probe does not perceive delay greater than  $\Delta_i$  as in Section III-D, the exit gateway returns an unchoke indication to the entry gateway. Otherwise, the exit gateway simply drops the probe. Observe that implementing choking and unchoking as above automatically handles path failures—(1) no data packets are sent on the failed path, and (2) if an unchoke request does not reach the exit gateway, no unchoke indication is sent back and hence the path remains choked.

#### F. Sender Rate Control

Harp is oblivious to the congestion control algorithm running at the sender. This can lead to a mismatch between the congestion windows on the sender and the entry gateway. The mismatch can lead to packet losses that reduce the congestion window on the sender, thereby reducing the throughput achieved by the background transfer.

One way to overcome the mismatch is to ensure that the sender does not send more bytes than what the entry gateway’s congestion window can allow across all paths. To achieve this effect, the entry gateway rewrites the TCP header in the acknowledgments returning to the sender with a receiver window equal to the minimum of the window allowed by the receiver and the congestion window allowed by the entry gateway. To handle receiver window scaling used by most bulk transfer applications [25], we monitor SYN packets to check if end-hosts exchange the scaling option, and scale the receiver window that we rewrite in the acknowledgments accordingly.

#### G. Implementation

We have implemented a prototype of Harp using the Click router toolkit [21]. Click provides the necessary infrastructure to capture, modify, route, and transmit packets. Using Click’s features, we build Harp such that it can be run both at the user- and the kernel-levels. Although for a given end-to-end connection our mechanism is distributed across different transit nodes,

each transit node can simultaneously function as entry, relay or exit node for different end-to-end connections.

*Transit node initialization:* In our implementation of Harp, an end-host (sender or receiver) using the mechanism associates itself to a transit node as its entry/exit node. This association can be done explicitly by the end-host, or seamlessly by an administrator. The information about each  $\langle \text{end\_host}, \text{associated\_transit\_node} \rangle$  pair is exchanged among transit nodes to enable successful routing of a sender’s packets through the transit nodes to the receiver. The configuration is similar to a firewall table; the identity of the end-host can be a wild-card to associate a group of IP addresses and (optionally) TCP ports to a transit node.

*Initial path selection:* For each entry and exit node pair, our prototype selects paths through relays that differ in RTTs from the direct path by at most a factor of 1.5.

*Background transfer identification:* Our current prototype assumes that the administrator of the entry node specifies a configuration similar to a firewall table to identify background flows.

#### IV. EXPERIMENTAL EVALUATION

We validate the effectiveness of Harp using both ns simulations and evaluation of a prototype on the Planetlab [22] testbed. Each of these environments allows us to demonstrate different aspects of Harp effectively. Further, simulations helped determine the right parameter settings for our real system prototype. Throughout this section, unless specified, we use the parameter settings as shown in Table I.

##### A. ns Simulations

In this section, we demonstrate different aspects of Harp using several topologies (shown in Figure 4) and a random topology generated by GT-ITM [26]. Topology1 is designed to (1) demonstrate that Harp can utilize spare capacity on multiple paths without causing much interference to foreground transfers, and (2) perform sensitivity analysis to various system parameters. Topology2 is designed to demonstrate that Harp leads to better allocation of network resources among multiple background transfers. Topology3 is generated by GT-ITM [26], and is used to demonstrate the effectiveness of Harp in complex traffic scenarios. In the first two topologies, T1 to T5 represent the transit nodes, and S# and R# represent source and destination of bulk transfers respectively. Topology3 contains 100 nodes acting as routers, and an additional 20 nodes connected to randomly chosen routers and acting as gateways and relays. The link bandwidths between routers vary from 10Mbps to 100Mbps. Finally, several clients and servers are connected randomly to routers.

The clients and the servers use TCP-Reno as the transport protocol and generate foreground traffic. All routers (represented as shaded circles) perform droptail FIFO queuing. The bottleneck links on each path are represented by dotted lines.

*Improvement in completion times:* We first show that Harp speeds-up background transfers substantially compared to

Param	Value	Remarks
$\alpha$	0.03	MI parameter, chosen as a small value to ensure stability [17], but large enough to get good throughput.
$\beta$	0.125	MD parameter, same as DECbit scheme [27], scalableTCP [28] and VCP [29]
$\delta$	0.15	Threshold on packet delays above which exit gateway indicates congestion. Value borrowed from TCP-LP study [14].
$\rho$	1.5	Reorder delay factor that determines how long an exit gateway delays packets.

Table I. Parameters used for simulations and prototype experiments.

unipath protocols, while causing minimum interference to foreground transfers. We use Topology1 and replay a web proxy trace between server and clients on both paths. The trace has 32 clients and contains requests to files of sizes varying from 150 Bytes to 1.7 MB, with an average size of 7.8 KB.

Figure 5(a) shows the completion time for a 100MB background transfer between S1 and R1 with Harp (Harp-multipath), with a unipath protocol (Harp-singlepath), and with TCP-Reno (Reno-singlepath) as we activate different number of clients on each path. The graph shows that, as the number of foreground clients increase, the completion time for the transfer increases in all cases, as expected. However, using multiple paths reduces the background completion time by a factor of two. Observe that Harp-singlepath achieves lower completion time than Reno-singlepath because Reno waits till a loss occurs and halves its congestion window every time the loss occurs. In contrast, Harp-singlepath avoids losses by backing off early (by a smaller proportion defined by  $\beta$ ).

To demonstrate that Harp reduces interference significantly compared to a non-background protocol such as Reno, we measure the *stretch* observed by the client requests. We define stretch for each request as the ratio of completion time for client requests when background transfers are active and completion time for the same requests when the transfers are not active. We compare the 90-th percentile stretch observed when the background transfers use one of Harp and TCP-Reno. Figure 5(b) shows that the stretch can be greater than 100% even with a single background transfer using Reno, and it rapidly increases to 500% with 29 transfers. With Harp, the stretch does not increase much with increasing background transfers, and is below 35% even with 29 transfers, thereby demonstrating that Harp transfers compete only for the spare capacity.

*Spare capacity utilization:* Figure 5(c) shows the effectiveness of Harp in quickly harvesting spare capacity left by a long running foreground transfer. The foreground transfer starts at  $t = 50s$  and ends at  $t = 200s$  and its sending rate is varied during this period as shown in the figure (by controlling its receiver window size). The graph shows that the throughput obtained by Harp closely reflects the spare capacity available with time.

*Fairness in resource allocation:* We first compare the variants (Independent, Joint and Adaptive) of multipath congestion control and show that the Adaptive variant chosen by Harp achieves fairness without compromising utilization. We consider Topology2 where we set up two background transfers: one between S1 and R1 that uses all four paths available (direct and through T2, T4 and T5), and the other between S2 and R2 that uses only one path (S2 to router to T3 to R2). Both transfers share

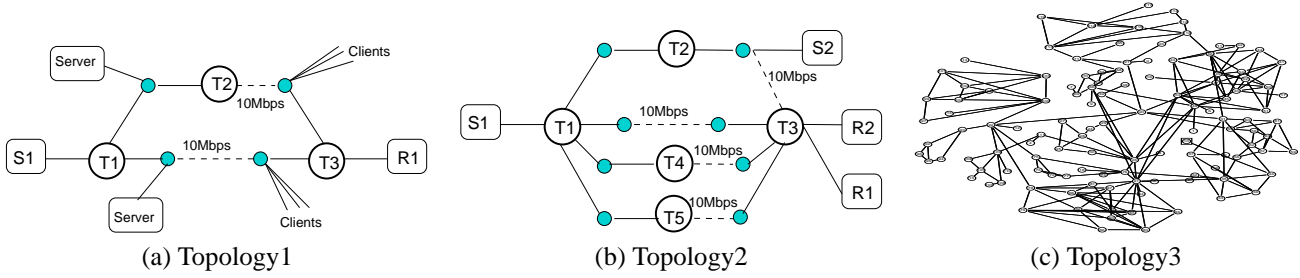


Fig. 4. Topologies for ns simulations.

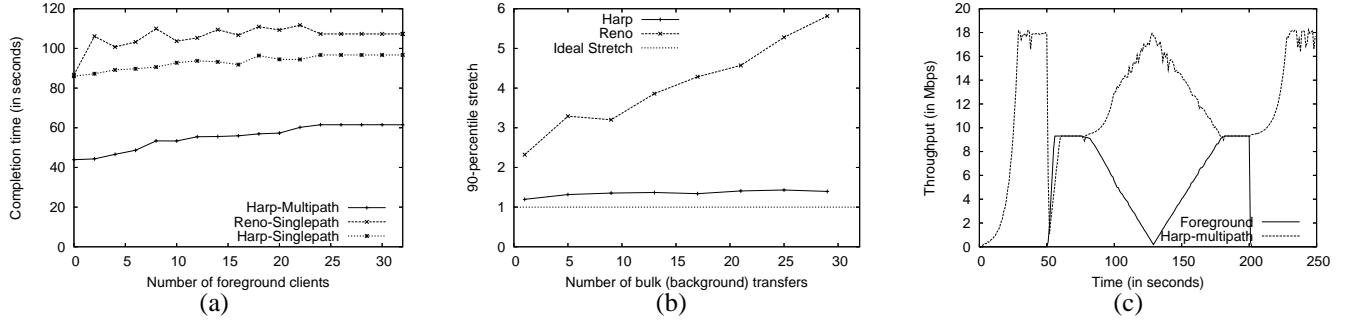


Fig. 5. (a) Spare capacity utilization with foreground web transfers, (a) Stretch caused to foreground requests, (c) Utilization with a long running foreground transfer

a bottleneck link. The transfer S2R2 starts 100 seconds after S1R1. Figure 6(a) shows the throughput achieved by both the transfers over time. Between 0 and 100s when S1R1 is the only active transfer, the throughput achieved when using Independent and Adaptive is about 20% higher than Joint. When S2R2 becomes active, both Adaptive and Joint yield throughput on path through T2 to ensure that S2R2 gets a better share of resources, whereas Independent does not yield the path thereby being significantly unfair to the transfer S2R2. Furthermore, Adaptive continues to utilize the other three paths more effectively than Joint between 100 and 250s. Note that in this experiment (and the others in this section) line 12 in Algorithm 1 never gets invoked because all transfers backoff based on delays and avoid losses; so the behavior of Joint is the same as the KV-controller [17]. These observations show the benefit of Adaptive in increasing utilization, while ensuring fair resource allocation when required.

Next, we use Topology2 and start multiple background transfers like S1R1 between T1 and T3, all of them utilizing all paths. Figure 6(b) plots Jain’s fairness index [30] for different number of transfers with time. The graph shows that the fairness index stays close to 1.0 with increasing number of transfers when all transfers share the same set of network resources.

*Utilization in a large topology:* We perform several experiments with Topology3 and different foreground and background traffic scenarios. For brevity, we present one result in Figure 6(c). We start 5 background transfers at different entry nodes (each lasting about 50 seconds) and 200 foreground transfers (each lasting about 1 second) at random times in a simulation run of 5 minutes. The graph shows the average throughput achieved by each foreground and background transfer. Harp, when alone (the bar Harp), reaps significant spare capacity in the network. The bars Harp/FG and Harp-S/FG show that (in the presence of foreground transfers) Harp with multiple paths reaps more

spare capacity than Harp with single path. The bars FG and FG/Harp represent foreground throughput with and without Harp transfers; Harp has little impact on foreground transfers.

#### A.1 Sensitivity Analysis

We study the sensitivity of Harp to resource parameters (such as varying bandwidths and link delays on paths), system parameters (such as reorder delay and  $\delta$ ), and design choices (such as choking and handling congestion window mismatch). For brevity, we discuss only a subset of the results. We perform all the experiments using Topology1.

*Heterogeneous link delays:* For this experiment, we vary the delay on the path T1,T2,T3 as a factor of the delay on path T1,T3; e.g., a link delay factor of 1.5 means T1,T2,T3 has 1.5 times longer one way delay than T1,T3. Figure 7(a) shows that as the variation in delays in the two paths increases, the completion times achieved by a multipath transfer (when alone and when competing with foreground clients replaying the proxy trace) increases. This graph hints that paths with low variation should be chosen in order to reap maximum benefits of multipath; further, paths with large variation can even be detrimental.

*Effect of reordering delay  $\rho$ :* Figure 7(b) shows the effect on completion times of maximum reordering delay  $\rho$  (Section III-C) applied to packets at the exit gateway. Small values of  $\rho$  result in greater number of packets getting reordered that make the receiver send duplicate acknowledgments and *falsely* indicate congestion to the sender. The sender backs off on such indications and thereby takes longer to complete the transfer. As  $\rho$  increases, the completion times decrease and flattens after certain value of  $\rho$ . Observe that in presence of foreground transfers, longer reordering delays are required to achieve maximum throughput. Using this graph, we pick  $\rho=1.5$  for our prototype.

*Benefit of choking:* To demonstrate the benefit of choking, we

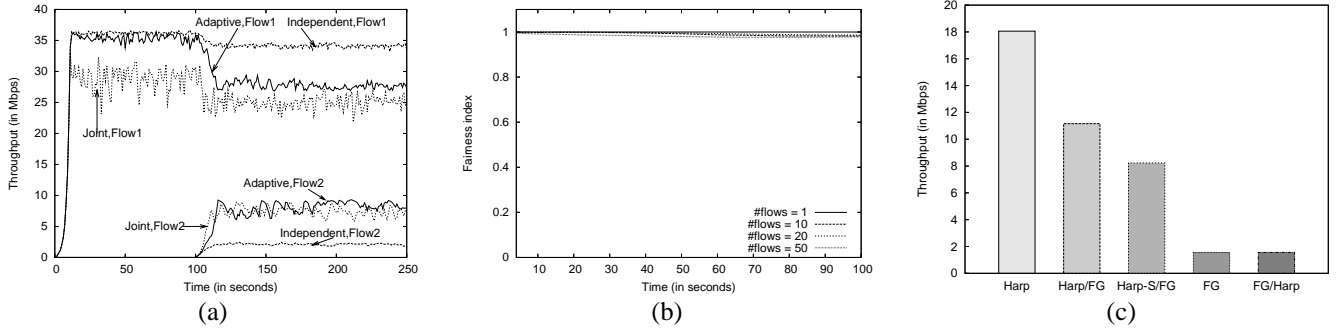


Fig. 6. (a) Resource allocation with different multipath congestion control algorithms, (b) Fairness among multiple multipath transfers, (c) Spare capacity utilization in a large topology

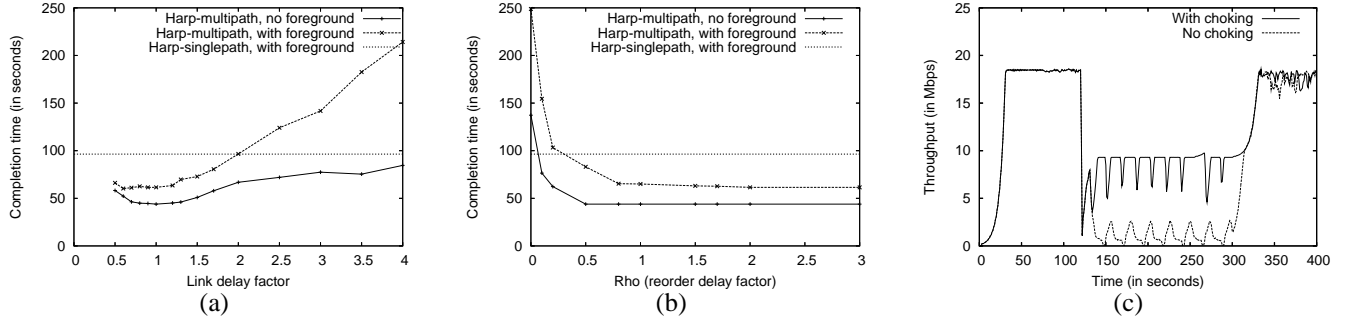


Fig. 7. (a) With increasing indirect path delay, (b) With increasing re-order delay, (c) With and without choking

plot in Figure 7(c) the throughput achieved with time by a multipath transfer between S1 and R1 (1) when using and (2) when not using choking. Between  $t=120s-320s$ , we start a long-running foreground transfer (not shown in the graph to avoid clutter) from server to a client through T2. The graph shows that when the foreground transfer starts, the bulk transfer backs off from the path T1,T2,T3. However, when choking is not used, the entry gateway sends one data packet per round-trip time to probe for available capacity. Since this packet will likely observe increased delay or even get dropped when the foreground transfer is active, packets sent on the direct path get delayed in the reorder queue and thereby reduce the throughput even on the direct path. In contrast, when choking is enabled, the entry gateway sends probe packets that do not contain any data, thereby maintaining high throughput on the direct path.

### B. Prototype Evaluation

In this section, we present an evaluation of Harp over Planetlab [22] to demonstrate its effectiveness in an Internet setting.

*Experimental Setup:* We setup an entry gateway and an exit gateway at Rutgers University and UMASS Amherst, and place a sender and a receiver on the same LAN as the entry and exit gateways respectively to avoid the access link of the sender and the receiver from becoming the bottleneck (see Figure 8(a)). We use several Planetlab nodes shown in Table II to act as relay nodes to provide path diversity between the entry and exit gateways. We divide the Planetlab nodes into three sets and repeat each experiment on all three sets. The right column in the table shows the range of minimum one-way delays between the entry and exit gateways through the relay nodes. We use Iperf [31] to

	Relay nodes	One-way Delays (ms)
Set 1	Columbia, Dartmouth, MIT	24-29
Set 2	NYU, Purdue UIUC	26-38
Set 3	Cornell, Princeton, RPI	23-32

Table II. Relay nodes, and one-way delays between entry and exit gateways through the relays.

generate foreground and background transfers of size 100MB.

*Results:* Figure 8(b) shows the completion times observed by the transfer when using one, two and three paths for each of the three sets. The results are an average of 10 runs. The solid bars show the average completion time. The figure shows that by using two additional paths, Harp decreases the completion time of a background transfer in all three sets by a factor of two.

To understand the impact of using multiple paths on packet reordering, we plot in Figure 8(c) the CDF of the time spent by the packets in the reorder queue at the exit gateway when using two and three paths. The graph shows that more than 90% of the packets spend less than 20ms in the queue, demonstrating that a small queue is sufficient to avoid most reordering. Also, less than 1% of the packets suffer a timeout and cause packet reordering.

To demonstrate the non-interference and fairness properties of Harp in simple realistic settings, we consider two sender receiver pairs S1R1 and S2R2 sharing links similar to Figure 3; S1R1 has access to two paths (through MIT and Dartmouth) and S2R2 has one path (through MIT) that is common with S1R1. For Figure 9(a), we setup one background transfer between S1R1 that uses the single path through MIT, and one foreground transfer between S2R2. We consider three scenarios—foreground alone, background alone, and foreground with back-

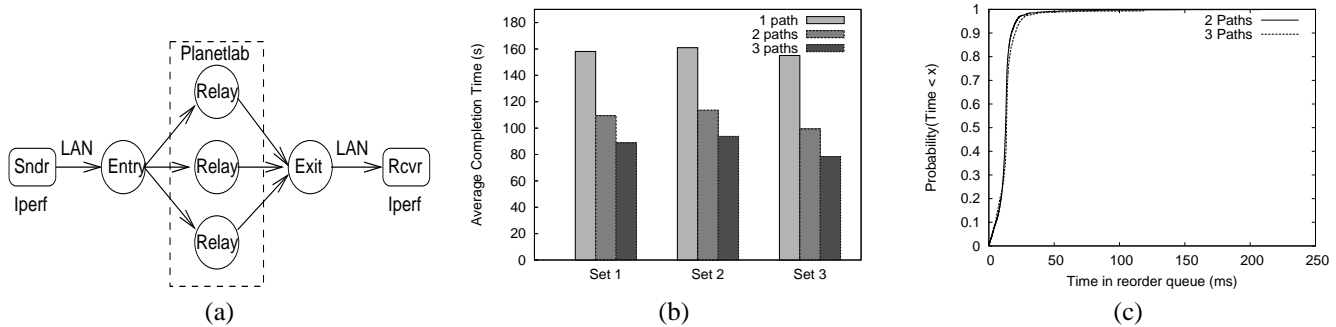


Fig. 8. (a) Wide-area experimental setup, (b) Through achieved by Harp with multiple paths, and (c) Reorder delay CDF at the exit gateway.

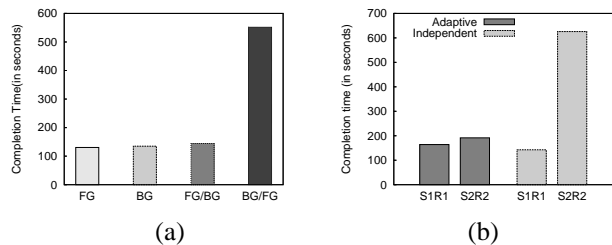


Fig. 9. (a) Non-interference, (b) Fairness

ground. We make the following observations: (1) When alone, background completion time (second bar) is close to the foreground completion time (first bar), thereby demonstrating that background reaps spare capacity effectively, and (2) foreground completion time with background (third bar) is close to foreground alone, thereby demonstrating that background effectively backs-off to make way for foreground; background’s completion time increases substantially as a result (fourth bar). For Figure 9(b), we setup a background transfer between S1R1 that uses both the available paths and a background transfer between S2R2, and examine two variants of multipath congestion control. The graph shows that S2R2 observes lower completion time with Adaptive than with Independent, thereby demonstrating that Adaptive leads to a fairer allocation of resources.

## V. RELATED WORK

This work bridges ideas from background transport protocols, multipath routing and congestion control, and overlay-based network services.

The primary goal of background transport is to be noninterfering with foreground transfers while providing practical levels of utilization of spare capacity. Noninterference is conceptually easy to provide using differentiated services, or Diffserv [12]. Once deployed, heterogeneous pricing can incentivize users to truthfully mark their traffic as foreground or background. However, it is unclear if and when router support to assist transport protocols will become available.

A more pragmatic effort in recent times has led to the development of end-host based protocols such as TCP Nice [13], TCP-LP [14], and Key et al. [15] that emulate low-priority transport without support from routers. The delay-based back-off in Harp’s controller is similar in spirit to Nice or LP. Our previous experience with Nice suggests that end-host based protocols are imperfect due to delayed feedback and can significantly

degrade utilization to maintain noninterference, e.g., in the presence of a competing Web workload, Nice can leave up to 50% of spare capacity on the path unused. Furthermore, even a perfect background protocol will find few takers without a suitable incentive mechanism for users to “nice” transfers.

Router-assisted as well as end-host background protocols fundamentally lack a mechanism for network-wide scheduling of resources and can arbitrarily delay background transfers along a path. Various studies have shown significant levels of spare capacity [13, 32], load imbalance [32], and background traffic in the Internet. Harp’s support for multipath background transport exploits this state of affairs to ensure quick completion of both foreground and background transfers.

Application-level approaches for background transport include (i) coarse-grain scheduling such as running backup [33] and data prefetching [34] applications during off-peak hours, (ii) rate limiting for applications such as Web prefetching [35, 36], (iii) adaptive rate control such as in Microsoft XP’s Background Intelligent Transfer Service [37] etc. These require tedious manual effort to configure parameters to reduce interference. Overlay quality-of-service approaches such as OverQoS [38] and MPAT [39] can provide limited noninterference among the set of flows they manage. In contrast, Harp provides a network-wide abstraction of low priority service. Harp’s relay nodes are simple and provide only forwarding functionality, while its gateway-based congestion control feature exists primarily to ease deployment at access points. Broadband users can directly use Harp without a gateway.

Multipath routing has seen a long history of work since Maxemchuk’s seminal work on dispersity routing [40] to improve throughput and resilience to path failures or packet losses. Striping or inverse multiplexing [41, 42, 43, 44] provides link level mechanisms for splitting input flows among multiple links to increase throughput. More recently, multipath techniques have been proposed in the context of overlay routing or multihomed clients [45, 46, 47]. Katabi et al. [48] and Elwalid et al. [49], propose splitting aggregate traffic flows along multiple paths to achieve load balancing and stability in the context of intradomain traffic engineering. In contrast, Harp uses multiple paths for scheduling transfers with heterogeneous requirements.

Lee et al. [50] and Zhang et al. [23] propose modifications to TCP’s fast retransmit or delayed acknowledgment schemes to tolerate packet reordering across heterogeneous paths. Harp interoperates with legacy TCP implementations by interposing

a network-layer connection proxy at a gateway or end-host.

Multipath transport mechanisms pTCP [51] and RMTCP [52] do not consider fairness when two or more paths for a single flow share the same bottleneck link. mTCP [23] by Zhang et al. explicitly identifies paths with shared bottlenecks and treats them as one flow. Thus, mTCP is TCP-fair on each path. However, it does not provide fairness of network-wide resource allocation, as users having access to more paths obtain an unfair share of network resources.

In contrast, Multipath TCP [18] and KV [17] are based on a utility-theoretic framework for network-wide resource allocation (see Kelly et al. [53] for a seminal paper), and systematically address fairness. These controllers maintain a nontrivial sending rate on each available path so as to maximize global fairness. Furthermore, these controllers implicitly behave like a single flow on paths with shared congestion. Our experimental effort to build a practical system for exploiting path diversity led us to address a limitation of the above controllers.

## VI. CONCLUSIONS

In this paper, we present the design and implementation of a network architecture, Harp, to allocate network resources more in line with human needs in a background-dominated Internet. Harp is based upon a multipath background congestion control algorithm that moves background traffic out of the way of foreground to ensure quick completion for both kinds of transfers. Compared to end-host protocols, this approach improves utilization and fairness of resource allocation. Our gateway-based implementation obviates modification to, or cooperation from, end-hosts. Our extensive evaluation using simulations and a prototype deployed on Planetlab demonstrates significant reduction in human wait-time.

## REFERENCES

- [1] "Kazaa." <http://www.kazaa.com>.
- [2] B. Cohen, "Incentives Build Robustness in BitTorrent," in <http://www.bittorrent.com/bittorrentecon.pdf>.
- [3] D. Kostić, A. Rodriguez, J. Albrecht, and A. Vahdat, "Bullet: high bandwidth data dissemination using an overlay mesh," in *SOSP*, 2003.
- [4] L. Cherkasova and J. Lee, "FastReplica: Efficient Large File Distribution within Content Delivery Networks," in *Proc. of USITS*, 2003.
- [5] R. Kokku, P. Yalagandula, A. Venkataramani, and M. Dahlin, "NPS: A Non-interfering Deployable Web Prefetching System," in *USITS*, 2003.
- [6] D. Rosenberg, "Prefetching Content for Increased Performance." [http://devedge-temp.mozilla.org/viewsource/2003/link-prefetching/index\\_en.html](http://devedge-temp.mozilla.org/viewsource/2003/link-prefetching/index_en.html), 2003.
- [7] "Imx series ip set top box." [http://www.matrixstream.com/IMX\\_1000.cfm](http://www.matrixstream.com/IMX_1000.cfm).
- [8] "LiveVault Corporation." <http://www.livevault.com/>.
- [9] "Remote backup systems." <http://remote-backup.com/rbackup/>.
- [10] "Microsoft Windows Update." <http://windowsupdate.microsoft.com/>.
- [11] CacheLogic Home Page, "Advanced Solutions for P2P Networks Home Page." <http://cachelogic.com/research/p2p2005.php>.
- [12] S. Blake et al., "An architecture for differentiated services," in *RFC 2475*, 1998.
- [13] A. Venkatramani, R. Kokku, and M. Dahlin, "TCP-Nice: A Mechanism for Background Transfers," in *Proc. of OSDI*, 2002.
- [14] A. Kuzmanovic and E. W. Knightly, "TCP-LP: A Distributed Algorithm for Low Priority Data Transfer," in *Proc. of INFOCOM 2003*, 2003.
- [15] P. Key, L. Massoulié, and B. Wang, "Emulating low-priority transport at the application layer: a background transfer service," *SIGMETRICS Perform. Eval. Rev.*, vol. 32, no. 1, 2004.
- [16] L. Peterson, S. Shenker, and J. Turner, "Overcoming the Internet Impasse through Virtualization," in *HOTNETS III*, 2004.
- [17] F. Kelly and T. Voice, "Stability of end-to-end Algorithms for Joint Routing and Rate Control," *SIGCOMM Comput. Commun. Rev.*, vol. 35, no. 2, pp. 5–12, 2005.
- [18] H. Han et al., "Multi-path TCP: A Joint Congestion Control and Routing Scheme to Exploit Path Diversity in the Internet," in *IMA Workshop on Measurements and Modeling of the Internet*, 2004.
- [19] "Akamai, Inc. Home Page." [www.akamai.com](http://www.akamai.com).
- [20] "Google Enterprise Services." <http://www.google.com/enterprise>.
- [21] E. Kohler et al., "The Click Modular Router," *ACM Transactions on Computer Systems*, vol. 18, August 2000.
- [22] "PlanetLab." <http://www.planet-lab.org>.
- [23] M. Zhang et al., "A Transport Layer Approach for Improving End-to-End Performance and Robustness Using Redundant Paths," in *Proc. of the USENIX 2004*.
- [24] R. Kokku, A. Bohra, S. Ganguly, and A. Venkataramani, "A Multipath Background Network Architecture." NEC Labs Tech. report. <http://www.nec-labs.com/~ravik/harp-tech.pdf>, 2006.
- [25] W. R. Stevens, *TCP/IP Illustrated, Volume 1: The Protocols*. 1994.
- [26] "Modeling Topology of Large Internetworks." <http://www-static.cc.gatech.edu/projects/gtitm/>.
- [27] K. K. Ramakrishnan and R. Jain, "A binary feedback scheme for congestion avoidance in computer networks with a connectionless network layer," in *SIGCOMM*, 1988.
- [28] T. Kelly, *Engineering flow controls for the Internet*. PhD thesis, University of Cambridge, 2004.
- [29] Y. Xia, L. Subramanian, I. Stoica, and S. Kalyanaram, "One more bit is enough," *SIGCOMM Comput. Commun. Rev.*, vol. 35, no. 4, 2005.
- [30] D. Chiu and R. Jain, "Analysis of the increase/decrease algorithms for congestion avoidance in computer networks," *Computer Networks*, 1989.
- [31] "Iperf, The TCP/UDP Bandwidth Measurement Tool." <http://dast.nlanr.net/Projects/Iperf/>.
- [32] A. Akella, S. Seshan, and A. Shaikh, "An empirical evaluation of wide-area internet bottlenecks," in *IMC*, 2003.
- [33] C. Maltzahn, K. Richardson, D. Grunwald, and J. Martin, "On bandwidth smoothing," in *4th International Web Caching Workshop*, 1999.
- [34] S. Dykes and K. A. Robbins, "A Viability Analysis of cooperative proxy caching," in *INFOCOM*, 2001.
- [35] N. Spring et al., "Receiver based management of low bandwidth access links," in *Proc of IEEE Infocom*, 2000.
- [36] M. Crovella and P. Barford, "The network effects of prefetching," in *INFOCOM*, 1998.
- [37] "Windows XP Background Intelligent Transfer Service." <http://www.microsoft.com/windowsserver2003/techinfo/overview/bits.msp>.
- [38] L. Subramanian, I. Stoica, H. Balakrishnan, and R. H. Katz, "Overqos: offering internet qos using overlays," *SIGCOMM Comput. Commun. Rev.*, vol. 33, no. 1, 2003.
- [39] M. Singh, P. Pradhan, and P. Francis, "MPAT: Aggregate TCP Congestion Management as a Building Block for Internet QoS," in *ICNP*, 2004.
- [40] N. F. Maxemchuk, *Dispersy routing in store-and-forward networks*. PhD thesis, Univ. Pennsylvania, Philadelphia, 1975.
- [41] H. Adishesu, G. Parulkar, and G. Varghese, "A reliable and scalable striping protocol," in *ACM SIGCOMM*, 1996.
- [42] J. Duncanson, "Inverse multiplexing," in *IEEE Communications Magazine*, vol. 32, 1994.
- [43] A. Qureshi and J. Gutttag, "Horde: separating network striping policy from mechanism," in *ACM MobiSys*, 2005.
- [44] K.-H. Kim and K. G. Shin, "Improving TCP Performance over Wireless Networks with Collaborative Multi-homed Mobile Hosts," in *ACM MobiSys*, 2005.
- [45] D. Andersen, A. Snoeren, and H. Balakrishnan, "Best-Path vs. Multi-Path Overlay Routing," in *Proc of IMC*, 2003.
- [46] A. Akella and J. Pang and B. Maggs and S. Seshan and A. Shaikh, "A Comparison of Overlay Routing and Multihoming Route Control," in *Proc of ACM SIGCOMM*, 2004.
- [47] A. Sen et al., "On Multipath Routing with Transit Hubs," in *Proc. of Networking*, May 2005.
- [48] S. Kandula, D. Katabi, B. Davie, and A. Charny, "Walking the Tightrope: Responsive Yet Stable Traffic Engineering," in *SIGCOMM*, 2005.
- [49] A. Elwalid, C. Jin, S. H. Low, and I. Widjaja, "MATE: MPLS adaptive traffic engineering," in *INFOCOM*, 2001.
- [50] Y. Lee, I. Park, and Y. Choi, "Improving TCP Performance in Multipath Packet Forwarding Networks," in *Journal of Communication and Networks*, vol. 4(2), Jun 2005.
- [51] H.-Y. Hsieh and R. Sivakumar, "A transport layer approach for achieving aggregate bandwidths on multi-homed mobile hosts," in *MobiCom*, 2002.
- [52] L. Magalhaes and R. Kravets, "Transport level mechanisms for bandwidth aggregation on mobile hosts," in *Proc. of ICNP*, Nov 2001.
- [53] F. Kelly, A. Maulloo, and D. Tan, "Rate control in communication networks: shadow prices, proportional fairness and stability," in *Journal of the Operational Research Society*, vol. 49, 1998.