

# Misbehaving TCP Receivers Can Cause Internet-Wide Congestion Collapse \*

Rob Sherwood  
Department of Computer  
Science  
University of Maryland,  
College Park  
capveg@cs.umd.edu

Bobby Bhattacharjee  
Department of Computer  
Science  
University of Maryland,  
College Park  
bobby@cs.umd.edu

Ryan Braud  
Department of Computer  
Science and Engineering  
University of California, San  
Diego  
rbraud@cs.ucsd.edu

## ABSTRACT

An *optimistic* acknowledgment (opt-ack) is an acknowledgment sent by a misbehaving client for a data segment that it has not received. Whereas previous work has focused on opt-ack as a means to greedily improve end-to-end performance, we study opt-ack exclusively as a denial of service attack. Specifically, an attacker sends optimistic acknowledgments to many victims in parallel, thereby amplifying its effective bandwidth by a factor of 30 million (worst case). Thus, even a relatively modest attacker can totally saturate the paths from many victims back to the attacker. Worse, a distributed network of compromised machines (“zombies”) attacking in parallel can exploit over-provisioning in the Internet to bring about wide-spread, sustained congestion collapse.

We implement this attack both in simulation and in a wide-area network, and show its severity both in terms of number of packets and total traffic generated. We engineer and implement a novel solution that does not require client or network modifications allowing for practical deployment. Additionally, we demonstrate the solution’s efficiency on a real network.

## Categories and Subject Descriptors

D.4.6 [Security and Protection]: Invasive Software

## General Terms

Security

## Keywords

Congestion Control, Distributed denial of service, Transmission control protocol

\*This work was supported by a grant (ANI 0092806) from the National Science Foundation.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CCS’05, November 7–11, 2005, Alexandria, Virginia, USA.  
Copyright 2005 ACM 1-59593-226-7/05/0011 ...\$5.00.

## 1. INTRODUCTION

Savage et al. [21] present three techniques by which a misbehaving TCP receiver can manipulate the sender into providing better service at the cost of fairness to other nodes. With one such technique, optimistic acknowledgment (“opt-ack”), the receiver deceives the sender by sending acknowledgments (ACKs) for data segments before they have actually been received. In effect, the connection’s round trip time is reduced and the total throughput increased. Savage et al. observe that a misbehaving receiver could use opt-ack to conceal data losses, thus improving end-to-end performance at the cost of data integrity. They further suggest that opt-ack could potentially be used for denial of service, but do not investigate this further.

In this paper, we consider a receiver whose *sole interest* is exploiting opt-ack to mount a distributed denial of service (DoS) attack against not just individual machines, but *entire networks*. In this paper, we:

1. Demonstrate a previously unrealized and significant danger from the opt-ack attack (one attacker, many victims) through analysis (Section 2.2) and both simulated and real world experiments.
2. Survey prevention techniques and present a novel, efficient, and *incrementally deployable* solution (Section 4.2) based on skipped segments, whereas previous solutions ignored practical deployment concerns.
3. Argue that the distributed opt-ack attack (many attackers, many victims) has potential to bring about sustained congestion collapse across large sections of the Internet, thus necessitating immediate action.

### 1.1 An Attack Based on Positive Feedback

Two significant components of transport protocols are the flow and congestion control algorithms. These algorithms, by necessity, rely on remote feedback to determine the rate at which packets should be sent. This feedback can come directly from the network [18, 11] or, more typically, from end hosts in the form of positive or negative acknowledgments. These algorithms implicitly assume that the remote entity generates correct feedback. This is typically a safe assumption because incorrect feedback rapidly deteriorates end-to-end performance [8]. *However, an attacker who does not care about data integrity could violate this assumption to induce the sender into injecting many packets into the network.* While not all of these packets may arrive at the receiver, they do serve to congest the sender’s network and saturate the path from the sender to the receiver.

Because acknowledgment packets are relatively small (40 bytes), it is trivial for an attacker to target hundreds and even thousands of victims in parallel. In effect, not only are each victims' access links saturated, but, due to over-provisioning, higher bandwidth links in the upstream ISPs begin to suffer congestion collapse in aggregate as well. In Section 2.4, we argue that sufficiently many attackers can overwhelm backbone links in the core of the Internet, causing wide-area sustained congestion collapse.

## 2. ATTACK ANALYSIS

In this section we describe pseudo-code for the attack, attack variants, and the details of the distributed version of the opt-ack attack. In Section 3, we present the observations we made in implementing the attack and techniques for mitigating practical concerns.

---

### Algorithm 1 –Attack( $\{v_1 \dots v_n\}, mss, wscale$ )

---

```

1: maxwindow ← 65535 × 2wscale
2: n ← |{v1, ..., vn}|
3: for i ← 1 .. n do
4:   connect(mss, wscale) to vi, get isni
5:   acki ← isni + 1; wi ← mss
6: end for
7: for i ← 1 .. n do
8:   send vi data request { http get, ftp fetch, etc. ... }
9: end for
10: while true do
11:   for i ← 1 .. n do
12:     acki ← acki + wi
13:     send ACK for acki to vi { entire window }
14:     if wi < maxwindow then
15:       wi ← wi + mss
16:     end if
17:   end for
18: end while

```

---

### 2.1 The Opt-Ack Attack

Algorithm 1 shows how a single attacker can target many victims at once. Typically, the attacker would employ a compromised machine (a “zombie” [23]) rather than launch the attack directly.<sup>1</sup> Consider a set of victims,  $v_1 \dots v_n$ , that serve files of various sizes. The attacker connects to each victim, then sends an application level request, e.g., an HTTP GET. The attacker then starts to acknowledge data segments *regardless of whether they arrived or not* (Figure 1). This causes the victim to saturate its local links by responding faster and faster to the attackers opt-acks. To sustain the attack, the attacker repeatedly asks for the same files or iterates through a number of files.

The crux of the attack is that the attacker must produce a seemingly valid sequence of ACKs. For an ACK to be considered valid, it must not arrive before the victim has sent the corresponding packet. Thus, the attacker must estimate which packets are sent and when, based only on the stream of ACKs the attacker has already sent. At first this might seem a difficult challenge, but the victim's behavior on receiving an ACK is exactly prescribed by the TCP congestion control algorithm! The attack takes three parameters: a

<sup>1</sup>This attack can also be mounted if the attacker is able to spoof TCP connections, either by being on the path between the victim and the spoofed address, or from guessing the initial sequence number, but we do not further consider it.

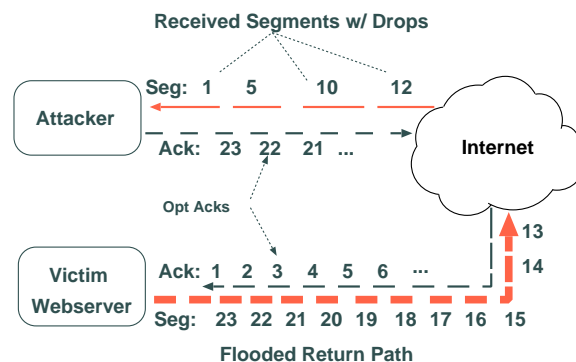


Figure 1: Opt-Ack Attack: Single Victim w/ Packet Loss (One of many victims)

list of  $n$  victims, the maximum segment size ( $mss$ ), and the window scaling ( $wscale$ ) factor. In the algorithm, the attacker keeps track of each victim's estimated window ( $w_i$ ) and sequence number to acknowledge ( $ack_i$ ). The upper bound of  $w_i$ ,  $maxwindow$ , is 65535 by default, but can be changed by the window scaling option (see Section 2.2). Note that the attacker can manipulate each victim's retransmission time out (RTO), because the RTO is a function of the round trip time, which is calculated by the ACK arrival rate. So, in other words, the attack can completely manipulate the victims in terms of how fast to send, how much to send, and when to time out.

There is a near arbitrary number of potential victims, given the pervasiveness of large files on the Internet. Any machine that is capable of streaming TCP data is a potential victim, including HTTP servers, FTP servers, content distribution networks (CDN), P2P file sharing peers (KaZaa[2], Gnutella[1]), NNTP servers, or even machines with the once common character generator (“chargen”) service.

The attack stream is difficult to distinguish from legitimate traffic. To an external observer that is sufficiently close to the victim, such as a network intrusion detection system (IDS), this stream is in theory indistinguishable from a completely valid high speed connection.<sup>2</sup> While it is common for IDSs to send out alerts if a large stream of packets enters the local network, the stream of ACKs from the attacker is comparatively small (see Section 2.2 for exact numbers). It is the stream of data *leaving* the network that is the problem.

Additionally, an attacker can further obscure the attack signature by sending acknowledgments to more victims less often, with the total amount of traffic generated staying constant. In other words, by generating less traffic per host and staying under the detection threshold, but increasing the total number of hosts *it is not locally obvious to the victims that they are participating in an DDoS attack*. As a result, short of globally coordination, it is difficult for victims to locally determine which if any of their data streams are malicious.

While Algorithm 1 works in theory, there are still challenges for the adversary to keep ACKs synchronized with the segments the victims actually send. We address these issues in the next section.

<sup>2</sup>Presumably, a monitoring system deployed closer to the attacker could detect the asynchrony between ACKs and data segments, but it is not practical to store per-flow state deep in the network.

## 2.2 Amplification

While it is not surprising that a victim can be induced to send large amounts of data into the network, the actual opt-ack amplification factor is truly alarming. The upper bound on the traffic induced across all victims from a single attacker is a function of four items: the number of victims ( $n$ ), and for each individual victim  $i$ , the rate at which ACKs arrive at each victim ( $\alpha_i$ ), the maximum segment size ( $mss_i$ ), and the size of the victim's congestion window ( $w_i$ ). Note that the attacker can use a single ACK to acknowledge an entire congestion window of packets. If we assume a standard TCP/IP 40 byte header with no options then the packet size is  $40 + mss_i$  bytes<sup>3</sup>. The rate of attack traffic in bytes/second is simply the sum across each victim of the product of the ACK arrival rate ( $\alpha_i$ ), the number of packets ( $\lfloor w_i/mss_i \rfloor$ ), and the size of each packet ( $40 + mss_i$ ). Given that the maximum window  $w_i$  is  $65535 \times 2^{wscale}$  and the total ACK arrival rate over all victims cannot exceed the attacker's bandwidth ( $\beta$ ) divided by the ACK size (40), we derive the total attack traffic  $T_{max}$  in bytes/second as:

$$T_{max} = \beta \times 65535 \times 2^{wscale} \times \frac{1}{mss} + \frac{1}{40} \quad (1)$$

As implied by (1), for typical Internet connections, i.e.,  $mss = 1460$  and  $wscale = 0$ , the attacker has an amplification factor of 1683. However, in the worst case, i.e., if  $mss = 88$  (the minimum  $mss$  for Linux and Windows XP[19]) and  $wscale = 14$ , then the amplification factor can reach as high as 32,554,441. In real world terms without window scaling, an attacker on a 56 Kilo-bits/s modem ( $\beta = 7000$  B/s) can generate 9,351,145 B/s or approximately 8.9MB/s of flooding summed across all victims. This value is more than the capacity of a T3 line, and close to the theoretical limit of a 100Mb Ethernet connection. With window scaling and small  $mss$ , a single attacker on a modem is capable of generating more traffic per second than the Slammer worm at its peak [14]. Obviously there are practical limits on the amplification including the maximum number of victims, their minimum bandwidth, and the time required to maximize the congestion window. Due to space considerations, an exploration of these constraints can be found in [22].

## 2.3 Lazy Opt-Ack

Lazy opt-ack is a variant of the standard opt-ack attack. As we will discuss in Section 3, the main difficulty in our implementation is in remaining synchronized with the sender's sequence number. The synchronization issue can be totally avoided if the attacker ACKs any segment that it actually receives, skipping missing segments. This lazy variant is malicious in that the attacker is effectively concealing any packet loss, thereby creating a flow that does not decrease its sending rate when faced with congestion (i.e., a non-responsive flow). Since the attacker is using the actual  $RTT$  to the victim, it generates less traffic than the attack described in Algorithm 1. However, it is well known [7] that in a congested network, a non-responsive flow can cause compliant flows to back off, creating a DoS. Note that the lazy variant is different from the standard attack in that it is impossible for the attacker to overrun the victim. This observation is precisely what makes many existing solutions insufficient. The skipped segments solution we provide in Section 4.2 protects against both the lazy and standard attacks.

<sup>3</sup>Additional header from the link layer may affect the packet size, but we do not consider it further.

## 2.4 Distributed Opt-Ack Attack

In this section, we consider the distributed case where *multiple attackers* run the opt-ack attack in parallel, trivially, and with devastating effect. The only coordination required is that each attacker chooses a different set of victims. Because a single attacker can solicit an overwhelming number of packets (as shown in Section 2.2) *a relatively small group of attackers can cause the Internet to suffer widespread and sustained congestion collapse.*

First, because opt-ack targets any TCP server, there are *millions* of potential victims on the Internet. Considering P2P file distribution networks alone, Kazaa and Gnutella have over 2 million [10, 9, 20] and 1.4 millions [12] users respectively that each host large multimedia files. While P2P nodes are typically low bandwidth home users, the popular content distributor Akamai runs over 14,000 [4] highly provisioned, geographically distributed servers.

It is not immediately clear how much traffic is necessary to adversely affect the wide-area Internet. One data point is the traffic generated from the Slammer/Sapphire worm. In [14], Moore et al. used sampling techniques to estimate the peak global worm traffic at approximately 80 million packets per second. At 404 bytes/packet, the worm generated approximately 31GB/s of global Internet traffic. Subsequent email exchanges by Internet operators [15] noted that many access links were at full capacity, and completely unusable. As noted in Section 2.2, it is theoretically possible for *a single attacker on a modem* to generate more than enough traffic to exceed this threshold using large  $wscale$  values. If using large  $wscale$  values were infeasible (for example, if packets containing the  $wscale$  option were firewalled), then *five attackers* on T3 connections with more typical TCP options, i.e.,  $mss = 1460$  and  $wscale = 0$ , would be sufficient to match the Slammer worm's traffic. If each attacker targeted sufficient number of victims, such that the load on no one victim was notably high, it would be difficult to locally distinguish malicious and valid data streams. So, unlike Slammer, there would be no clear local rule to apply to thwart the attack.

The traffic from the Slammer worm was not sufficient to push the core of the Internet into congestion collapse. Because of the inherent difficulty in modeling wide scale Internet phenomena, it is not clear how to estimate the number of opt-ack attackers required to induce such a collapse. However, a single attacker on a modem or a small number of other attackers can induce traffic loads equivalent to the Slammer worm. Recent studies[5] show that there exists networks of compromised machines ("botnets") with over 200,000 nodes. Since each of these nodes represents a possible attacker, a large distributed opt-ack attack could easily be catastrophic.

## 3. IMPLEMENTING OPT-ACK

The main challenge in implementing the attack is to accurately predict which segments the victim is sending and ensure that the corresponding ACKs arrive at the correct time. In Figure 1, the attacker injects ACKs into the network before the corresponding segments have even reached the attacker, so remaining synchronized with the victim can be non-trivial. Maintaining this synchronization of sequence numbers is crucial to the attack. If the attacker falls behind, i.e., it starts to acknowledge segments slower than they are sent, then the victim slows down, may time out, and the effect of the attack is reduced. Similarly, if the attacker gets ahead of the victim in the sequence space, i.e., the victim received ACKs for segments that are not yet sent, the victim ignores these ACKs and the stream stops making progress. We refer to this condition as *overrunning* the victim. Overruns can occur in three different ways: ACKs arriving too quickly, lost ACKs, and delays at the

server. Below, we describe a technique for the attacker to detect this condition and recover.

In accordance with RFC793 [3], Section 3.4, when the sender receives ACKs that are not in the window, it should not generate a RST, but instead an empty packet with the correct sequence number. One of the tenets of the Internet design philosophy is the robustness principle: “be conservative in what you send, and liberal in what you accept,” and it is this principle that opt-ack exploits.

There are many ways that an overrun condition may result, most common being the sending application stalls its output because it was preempted by another process. In general, there are a myriad of factors that affect the sender’s actual output rate, including: the victim’s load, application delay, the victim’s send buffer size, and the victim’s hardware buffer. However, these factors are mitigated when the number of victims is large. By sending ACKs to more victims, each individual victim receives ACKs less often. This provides more time for the victim to flush its buffers, place the sending application back into the run queue, etc.

It is worth noting that the implementation we developed is only a demonstration of the potential severity of opt-ack. It is by no means an optimal attack. There are a number of points where a more thorough attacker might be able to mount a more efficient attack. However, as we note in Section 5, the implementation is sufficiently devastating as to motivate immediate action. Below, we discuss further strategies to mitigate and recover from overrunning the victim.

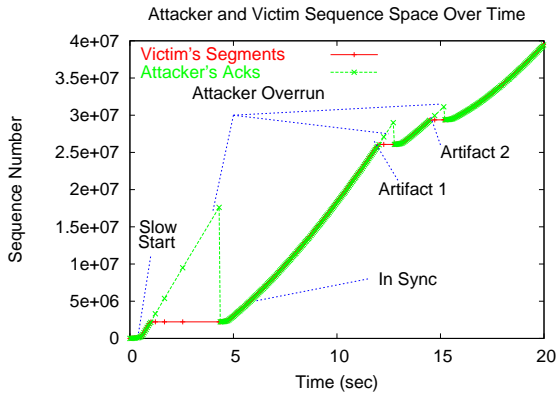


Figure 2: Attacker and Victim Sequence Space, Measured at Victim

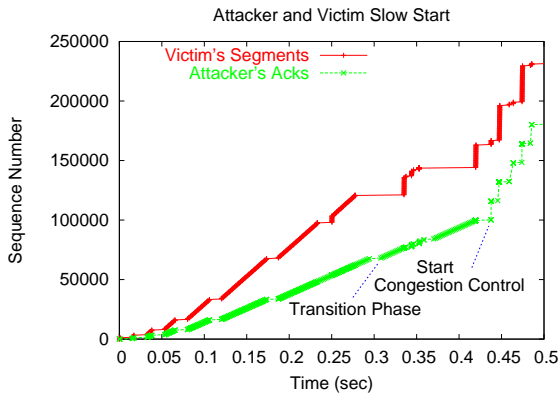


Figure 3: Detail: Attacker and Victim Slow Start

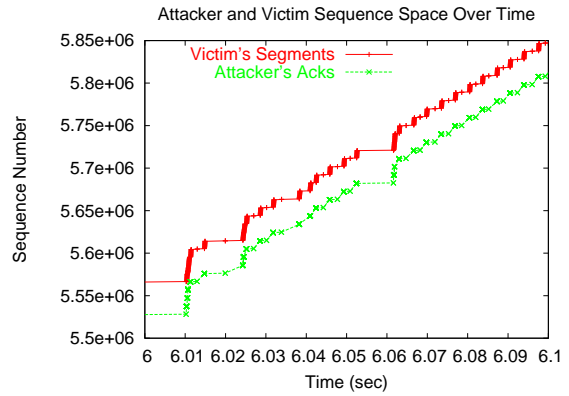


Figure 4: Detail: Attacker and Victim Synchronized

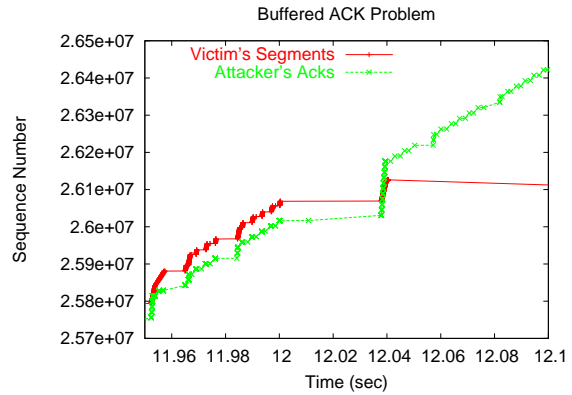
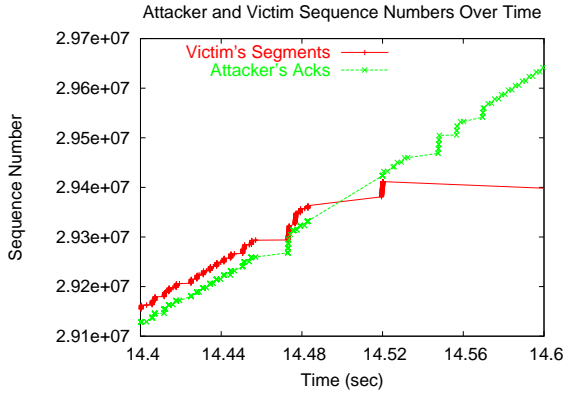


Figure 5: Artifact 1: Buffered ACKs

### 3.1 Recovery from Overruns

Compliant TCP streams are supposed to generate an empty segment upon receipt of an out of window ACK [3]. The attacker could use this empty segment to detect overruns, but this technique is unreliable because the incoming link is typically saturated, so the empty segment will be dropped before it reaches the attacker. Additionally, Linux’s TCP implementation does not send empty segments, instead ignores an out of window ACK ( Other OSes, specifically MacOS X 10.2 and Windows 2000, correctly generate the empty packet). However, note that in a stream that is making progress, i.e., not overrun, the sequence numbers of packets received increase monotonically (barring packet reordering). Upon a retransmission, or when an empty packet is received, the sequence number is less than or equal to the previous packet, breaking monotonicity. So, by monitoring the sequence numbers of packets actually received, the attacker can declare an overrun when the sequence numbers no longer increase. When an overrun is detected, the attacker can resume slow start on the last received packet, thus recovering from the overrun. This is an expensive process, as it potentially requires waiting on the order of at least one second [6] for the server to timeout.

Figure 2 shows the life cycle of an attack against a GNU/Linux 2.4.20 victim, across a wide area network, as measured at the victim. The “attacker” data points show the ACKs at the time the victim received them, and the “victim” data points show the segments being sent by the victim. Note that for the majority of the time the two lines are indistinguishable, i.e. the attacker is synchronized with the victim (Figure 4). However, on three occasions



**Figure 6: Artifact 2: Victim Delay and Buffered ACKs**

the attacker overruns the victim’s sequence number, and is forced to recover, as described above. The attacker blindly continues sending ACKs that are ignored, as the victim stops making progress in sending the stream (as demonstrated by the flat line). In the first overrun, the victim actually retransmits three times before the attacker recovered, because the retransmitted packets were also lost. However, in the next two overruns, the attacker recovered faster, each on the order of one second.

Recovery code must track the victim’s slowstart threshold ( $ssthresh_i$ ) in addition to the estimated congestion window ( $w_i$ ). The variable  $ssthresh_i$  is initialized to the maximum window size, is set to half  $w_i$  with every recovery, and grows with the congestion window, as prescribed by [24].

### 3.2 Victim’s Processing Time

One of the most difficult challenges in keeping the attacker synchronized is estimating the time spent for the victim to send the packets, which we call the processing time. Obviously, an attacker should not ACK segments faster than a victim is capable of generating them. If the attacker knows the victim’s processor speed, server load, operating system, and local bandwidth, it may be able to estimate the processing delay time. However, this information is difficult to determine, and underestimating the delay time leads to the attacker overrunning the server which causes reduced attack strength. To address this challenge, we introduce the *TargetBandwidth* variable. With this variable, we can derive the processing delay:

$$processing\ delay = \frac{\lfloor cwnd/mss \rfloor \times (40 + mss)}{TargetBandwidth}$$

The *TargetBandwidth* variable represents the rate of traffic the attacker is trying to induce the server to generate (in bytes/second). While the value of *TargetBandwidth* can be determined adaptively based on how often the attacker is forced to recover, for the purposes of the implementation code, we specify it as a runtime parameter. Additionally, many TCP implementations will locally invoke congestion control measures if the local hardware buffers become full, so it is important to pick the *TargetBandwidth* to be something feasible by the victim’s hardware.

The processing time of an idle server is significantly shorter than that of a busy server. This implies that an attacker needs to estimate a server’s load before attacking it. However, we noted that as the attacker’s flow rate increases, the other connections are forced to back off, which in turn decreases the processing time of the server. Thus, we introduce the concept of adaptive delay. By overesti-

imating the initial processing time and the delay between ACKs, i.e. sending ACKs slowly, and then progressively ramping up the ACK speed to the desired rate, third party streams are “pushed” out of the way with minimal overruns. How to do this effectively in an aggressive manner, without causing the attacker to overrun and restart, is an open question. However, in the implementation, we start arbitrarily at 10 times the estimated processing time, and then decrease down to the target processing time in steps of 500  $\mu s$  per window.

Another variable affecting the processing time is the coarse grained time slice in the victim’s scheduler. Periodically, the victim process is suspended for a number of time slices, which can cause a delay in sending if the kernel buffer is drained before the process can be rescheduled. An example of this is the second artifact (Figure 2, blown up as Figure 6), where the server actually pauses for 36 ms. Note, it is less obvious from Figure 6, but the server starts sending less than one millisecond before the buffered ACKs arrive. We do not have a technique to predict these delays, and rely on the recovery/restart mechanism.

### 3.3 Multiple ACKs Per Window and the Transition Phase

We noted that during congestion avoidance, the server rarely sent a full 64KB window, even when the congestion window would otherwise have allowed for it. The effect was that the number of segments in flight varied, and it became difficult for the attacker to ACK the correct number of segments. We speculate this is due to operating system buffering inefficiencies, and perhaps coarse grained time slices. Whatever the reason, we changed the attack algorithm to ACK half of the window at a time with the appropriate delay instead of the full window all at once. By ACKing half as much, twice as often, we were able to keep the amount of flooding high, reducing the chance the attacker overruns the victim’s sequence number. However, by sending twice as many ACKs, the attacker is restricted to half of the amplification listed in Section 2.2.

An additional effect of sending two ACKs per window is resistance to lost ACKs. The basic algorithm assumes that each ACK successfully reaches the victim, which is obviously not true on the general Internet. To maximize robustness, the implementation sends two ACKs offset by one *mss* from each other twice per window for a total of four ACKs per window. The effect here is two fold. First, the attacker can now lose three sequential ACKs in a row without overrunning the server. Second, with more ACKs the congestion window grows faster after recovery from overrun. Note that sending four ACKs per window is we reduce the expected amplification by a factor of four.

In development, it was difficult to track the exact state of the victim’s congestion window and  $ssthresh$ , especially after recovery. It was common for the attacker to stay correctly synchronized with the victim through slow start and then get out of sync immediately when moving to the congestion avoidance algorithm. While we speculate there are many factors that cause this behavior, i.e. unpredictable server load, and the timing involved in the congestion avoidance phase may need to be more accurate than the slow start phase, it simply became easier to work around it. Thus, we introduce a “transition” phase for the attacker between slow start and congestion avoidance (see Figure 3). In this transition phase, we ACK every expected packet in turn for the full window. Effectively, the transition phase allows for a larger margin of error in estimating the victim’s  $ssthresh$  variable. In practice, we ACK two full windows in the transition phase before transitioning to the full congestion avoidance portion of the attack.

### 3.4 The Attacker’s Local Bandwidth

Algorithm 1 does not take into account the attacker’s local bandwidth. Given a local bandwidth of  $\beta$  in bytes per second, the sum of all ACKs sent to all victims can be sent at at most  $\alpha = \beta/40$  bytes/second. At speeds faster than  $\alpha$ , and the ACKs get buffered or even dropped, which interferes with the timing of the attack. When ACKs are buffered (as shown in the first artifact of Figure 2, and Figure 5) they arrive at the victim all at once. The victim is not able to send fast enough to keep up with the sudden flood of ACKs and this creates an overrun. To fix this, we limit the rate of outgoing ACKs from the attacker as a function of the available local bandwidth, which is specified at runtime. The main effect of rate limiting the ACKs is to maintain even spacing when they arrive at the victim, despite network jitter and buffering. Additionally, as the number of victims increases, the difficulty from buffered ACKs decreases because the time between when the server receives ACKs increases.

## 4. DEFENDING AGAINST OPT-ACK

In this section, we present a simple framework for evaluating different defense mechanisms against the opt-ack attack, and evaluate potential solutions within that framework. Finally, we present one particular solution, randomly skipping segments, that efficiently and effectively defends against opt-ack. We also describe an implementation of randomly skipped segments in detail.

### 4.1 Solutions Overview

Any mechanism that defends against opt-ack should minimally possess the following qualities:

1. **Easy to Deploy** Due to the severity of the attack, any solution should be practically and quickly deployable in the global infrastructure. Minimally, the solution should allow incremental deployment, i.e., unmodified clients should be able to communicate with modified servers.
2. **Efficient** Compliant (i.e., non-attacking) TCP streams should suffer minimal penalty under the proposed solution. Also, low power embedded network devices do not have spare computational cycles or storage space. Because the problem is endemic to all implementations, the solution needs to be efficient on all devices that implement TCP.
3. **Robust** Any fix needs to defend against all variants (Section 2.3) of the opt-ack attack.
4. **Easy to Implement** This is a more pragmatic goal, leading from the observation that TCP and IP are pervasive, and run on a diverse range of devices. Any change in the TCP specification would affect hundreds (or thousands) of different implementations. As such, a simpler solution is more likely to be implemented.

In [22], we compare the costs and benefits of many defenses including secure nonces, ACK alignment, bandwidth caps, in network support, disallowing out of window ACKs, and random pauses. Table 1 is a summary of the defenses, and we present the most relevant of these solutions in detail below.

#### 4.1.1 Secure Nonces

One possible solution is to require that the client prove receipt of a segment by repeating an unguessable nonce. Assume each outgoing segment contains a random nonce which the corresponding ACK would have to return in order to be valid. Savage [21]

et al. improve on this solution with *cumulative* nonces. In their system, the response nonce is a function of all of the packets being acknowledged, i.e., a cumulative response, ensuring that the client actually received the packets it claims to acknowledge.

Unfortunately, cumulative nonces are not practically deployable. They requires both the client and server to be modified, preventing incremental deployment. If deployment was attempted, updated servers would be required to maintain backward compatibility with non-nonce enabled clients, until all client software was updated. As a result, updated servers would have to chose between being vulnerable to attack or compatibility with unmodified clients. Additionally, nonces require additional processing and storage for the sender. Calling a secure pseudo-random generator once per packet could prove expensive for devices with limited power and CPU resources, violating our efficiency goal.

To aid deployment, one could consider implementing nonces in existing, unmodified clients via the TCP timestamp option. The sender could replace high order bits of the timestamp with a random challenge, and any non-malicious client which implemented TCP timestamps would respond correctly with the challenge. If a client did not implement timestamps, the server could restrict throughput to something small, e.g, 4Kb/s. While this improves on the deployment of nonces, this solution still has problems. First, it loses the critical cumulative ACK property of Savage’s solution. That is, an acknowledgment for a set of packets does not necessarily imply that all packets in the set were received, which opens itself to the lazy opt-ack attack. Second, as we discuss in Section 4.1.2 below, bandwidth caps are not effective.

#### 4.1.2 Bandwidth Caps

The obvious solution to an attacker consuming too many resources, as is the case with the opt-ack attack, is to limit resource consumption. Conceivably, this could be done at the server with a per IP address bandwidth cap, but unfortunately this is not sufficient. First, any restriction on bandwidth can simply be over come by increasing the number of victims. Suppose for example, that each victim sets the policy that no client can use more than a fraction  $c \in (0, 1]$  of their bandwidth. Then the attacker need simply increase the number of victims by  $1/c$  to maintain the same total attack traffic. Further, bandwidth caps interfere with legitimately fast clients, violating our efficiency goal.

#### 4.1.3 Disallow Out of Window ACKs

A straightforward solution is to change the TCP specification to disallow out of window ACKs. Recall from Section 3 that our implementation runs the risk overrunning the victim. If a victim sent a reset, terminating the connection, upon receipt of an out of window ACK, the opt-ack attack would be mitigated. However, this is not a viable solution as this opens non-malicious connections to a new DoS attack. A malicious third party could inject a forged out of window ACK into a connection, causing a reset [25]. Because the ACK is out of window, there would be no need to guess the sequence space and window size. Also, compliant receivers can send out of window acknowledgments due to delays or packet reordering. For example, suppose ACKs for packets numbered 2 and 3 are sent but received in reverse order. The ACK for packet 3 would advance the window, and then the ACK for packet 2 would be and out of window ACK, causing a RST. Additionally, the lazy opt-ack attack is not prevented by disallowing out of window ACKs.



Solution	Efficient	Robust	Deployable	Simple	Change TCP Spec.
Cumulative Secure Nonces	yes	yes	no	yes	client & server
Secure Nonces w/ timestamps	yes	no	yes	yes	server only
ACK Alignment	yes	no	yes	yes	server only
Bandwidth Caps	no	no	yes	yes	no
Network Support	yes	yes	no	no	no
Random Pauses	no	no	yes	yes	server only
Skipped Segments	yes	yes	yes	yes	server only

Table 1: Summary of Defenses to Opt-Ack Attack

## 4.2 Proposed Solution: Randomly Skipped Segments

To defend against the opt-ack attack, we propose that the server randomly *skip* sending the current segment, and instead send the rest of the current window. Note that this is equivalent to locally, *intentionally* dropping the packet. A non-malicious client that actually gets all of the packets, save the skipped one, will start re-ACKing for the lost packet, thereby invoking the fast retransmit algorithm. However, an attacker, because it does not have a global view of the network, cannot tell where along the path a given packet was dropped, so it *cannot tell the difference between an intentionally dropped packet and a packet dropped in the network by congestion*. Thus, an attacker will ACK the skipped packet, alerting the server to the attack. Note that usually fast retransmission indicates network congestion, so the congestion window is correspondingly halved. However in this case, retransmission was not invoked due to congestion in the network, so the sender should not halve the congestion window/slow start threshold as it typically would. Given that most modern TCP stacks implement selective acknowledgments (SACK)[13], this solution is very efficient (see Section 5 for performance). *The only penalty applied to a conforming client is a single round trip time in delay.*

To determine how often to apply the skipped packet test, we maintain a counter of ACKs received. Once a threshold number of ACKs are received, the skip test is applied. It is important that the threshold be randomized, as the security of this system requires that the attack not predict which segment was skipped. However, there is an obvious trade off in where to make the skipped packet threshold. If it is too low, the server will lose efficiency from skipping packets too often. Setting the threshold too high allows the attacker to do more damage before being caught (see Section 5 for an exploration of this trade-off). Our solution is to chose the threshold uniformly at random over a configurable range of values.

This simple skipped segment solution meets all of our goals. It is efficient: compliant clients suffer only one round trip time in delay, the computational costs consist of keeping only an extra counter, and the storage costs are trivial (5 bytes per connection, described in [22]). The skipped packet solution is robust against the variations of the attack described in Section 2.3, because it inherently checks whether a client actually received the packets. This solution is a local computation, so it needs no additional coordination or infrastructure, i.e., the deployment requirements are met. Best of all, it is transparent to unmodified clients, allowing for incremental deployment. Due to space considerations, we discuss of our implementation of randomly skipped segments for the Linux kernel in [22].

Last, we must take care to insure that the randomly skipped segments solution test does not introduce a new DoS attack. For example, an attacker might maliciously inject ACKs into a benevolent client’s TCP stream, causing the server to believe the client

is performing a opt-ack attack. This attack is easily remedied if servers ignore out of window ACKs during the skipped segments test. Thus, attackers must guess the current sequence space and window size in order to correctly inject a malicious ACK. Further, in our randomly skipped segments implementation, when the servers detects an opt-ack attack, it simply reset the connection. Because any attacker that can guess the sequence space and window size can already reset a connection by injecting a RST packet, our solution introduces no new vulnerabilities into the protocol.

## 5. ATTACK EVALUATION

We evaluate the feasibility and effectiveness of the opt-ack attack in a series of simulated, local area, and wide area network experiments. In the first set of simulations, we determine the total amount of traffic induced by the opt-ack attacks. Next, we determine the effect of the attack on other (honest) clients trying to access the victim. We also present results for the amount of traffic (described in Section 2.2) our real world implementation actually achieves across a variety of platforms. Finally, in Section 5.3, we evaluate the efficiency of our skipped segment solution. Further experiments with number of packets sent, variable file sizes, and performance of the randomly skipped segments solution without selective acknowledgments can be found in [22].

### 5.1 Simulation Results

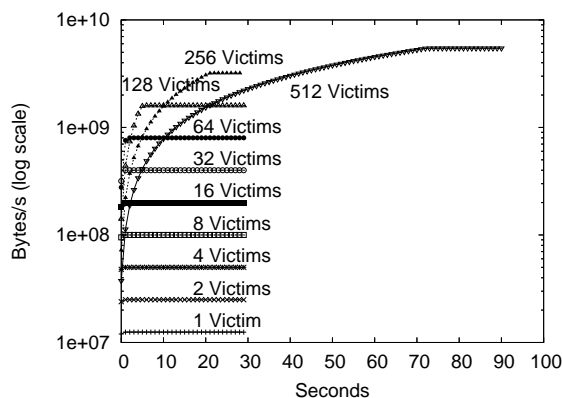


Figure 7: Maximum Traffic Induced Over Time; Attacker on T1 with  $mss=1460, wscale=4$

We have implemented the opt-ack attack in the popular packet level simulator ns2 and simulate the amount of traffic induced in various attack configurations. In each experiment, there is a single attacker and multiple victims connected in a star topology. Each victim has a link capacity of 100Mb/s, and all links have 10ms latency (the choice of delay is arbitrary because it does not affect the

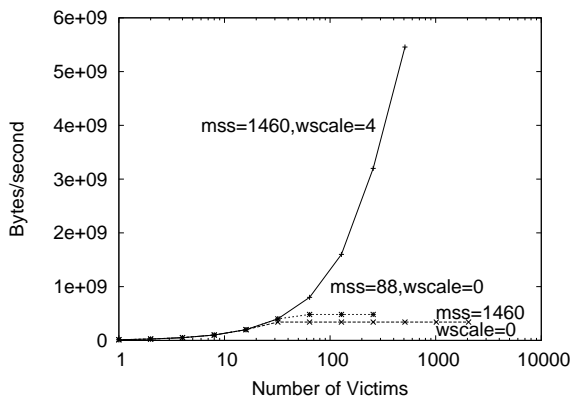


Figure 8: Maximum Traffic Induced By Number of Victims; Attacker on T1

attack). We vary the number of victims, and the  $mss$  and  $wscale$  of the connection. The attacker makes a TCP connection to each victim in turn, and only sends acknowledgments once all victims have been contacted. Victims are running the “Application/FTP” agent, which uses an infinite stream of data.

In Figure 7, we show the sum of the attack traffic generated over time with variable numbers of victims. In this experiment, the attacker is on a T1 (1.544Mbps) and uses connection parameters  $mss=1460$  and  $wscale=4$ . When the number of victims is less than 512, the amount of flooding is limited by the sum of the bandwidths of the victims. The amount of traffic doubles as the number of victims double until 512 victims. As the number of victim’s increases, the attack takes longer to achieve full effect. Both of these effects are discussed and analyzed in [22]. The case with 512 victims took 73 seconds to reach it peak attack rate, while all others did so in under 30 seconds. At 512 victims, the simulation achieves 99.9% of the traffic predicted by Equation 1.

As shown in Figure 7, once the attack’s maximum effect is reached, it can be sustained indefinitely. In Figure 8 we show the maximum traffic induced as we vary the number of victims,  $mss$  and  $wscale$  for bytes/second. As predicted by Section 2.2, attackers with a lower  $mss$  produce more traffic than one with a higher value. Likewise, an increased  $wscale$  has a dramatic increase in the total traffic generated. In [22], we also present results that measure the total number of packets/second sent.

Due to CPU and disk space limits, we were not able to simulate more than 512 victims for all parameters, or  $wscale$  values above 4, despite the fact that our simulation machine was a dual processor 2.4Ghz Athlon-64 with 16GB ram and 300GB in disk.

## 5.2 Real World Implementation

We implemented the attack in C and experimented on real machines in various network settings. We measure the actual bandwidth generated from a single victim running various popular operating systems. We did not experiment with multiple attackers or victims due to real world limitations of our test bed. Our experiments with a single attacker and single victim were sufficient to cause overwhelming traffic on our local networks. It would be irresponsible and potentially illegal to have tested the distributed attack on a wide-area test bed (e.g., PlanetLab[17]), and even our simple one attacker-one victim wide-area experiments caused network operators to block our experiments.<sup>4</sup>

<sup>4</sup>Incoming traffic to one author’s home DSL IP address was tem-

Experiment	Average (sec)	Dev.	Increase
No Attack	89.11	0.007	1
LAN Attack	1552.03	141.76	17.42
WAN Attack	779.93	139.32	8.75

Figure 9: Average Times with Deviations for a Non-malicious Client to Download a 100MB File

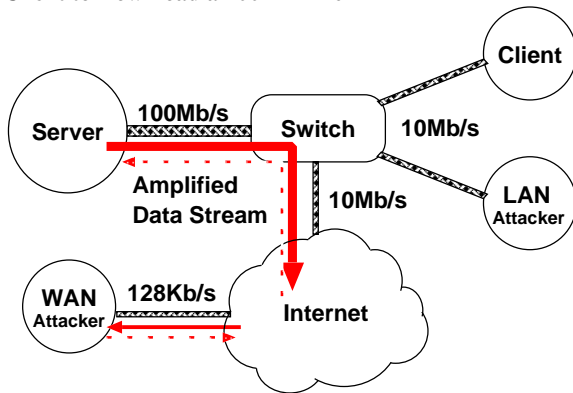


Figure 10: Topology for Experiments

### 5.2.1 Single Victim DoS Effect - LAN and WAN

This experiment measured the effect on a third party client’s efficiency in downloading a 100MB file from a single victim during various attack conditions. We repeated this experiment with no attacker, with an attacker on the local area network, and with an attacker across the Internet (see Figure 10). The local area attacker was a dual processor Pentium III running Linux with a 10Mb Ethernet card, while the WAN attacker was a 100Mhz Pentium running GNU/Linux on an asymmetric 608/128 Kb/s downstream/upstream residential DSL line. The latency on the WAN link varied over time, with a average RTT of 13.5ms.

A typical web server runs on a fast local area network, which connects to a slower wide area network. In order to emulate this bottleneck, and also to safeguard against saturation of our production Internet connection, we connected our test web server to the world via a 10Mb connection on a Cisco Catalyst 3550 switch. Furthermore, both LAN and WAN attackers were configured to use  $TargetBandwidth$  of  $10^9$  bytes/second, and  $\beta = 16000$  bytes/s as their local bandwidth setting (see Section 3 for description). The intuition is that the LAN and WAN attackers should be equally capable with respect to their available bandwidth, but the WAN attacker must compensate for more end-to-end jitter and delay. Each run used  $mss=536$  and  $wscale=0$ , i.e., typical values for Internet connections. Each experiment was repeated 10 times and the values averaged. The numbers were measured with a command line web client (similar to *wget*) specially instrumented to measure bandwidth at 10 ms intervals. We present the results from these experiments in Table 9. The “Increase” column refers to the increase in time relative to the “No Attack” baseline.

The effect of the attack is significant. The 100MB file takes on average 17.42 and 8.75 times longer to download under LAN and WAN attack, respectively. We believe that the time difference between the WAN and LAN attacks is due to the increased jitter of the wide area Internet, and the increased standard deviation in the

porarily blocked as a result of these experiments. This did not serve to stop the attack, as the outbound ACKs could still be sent. However, this served as evidence that we should cease the experiment.



OS	Avg. KB/s	Dev.	Amplification
Linux 2.4.24	3931.93	1102.38	251.6
Mac OS X	806.2	258.1	51.6
Solaris 5.8	3150.6	1301.1	201.6
Windows XP	640.62	378.85	41.0

**Table 2: Average kilobytes/s of Induced Flooding, Standard Deviation, and Amplification Factor of Attacker’s Bandwidth**

results supports this. This variability makes keeping synchronizing with the victim more difficult due to the buffered ACK problem, as described in Section 3. However, more advanced attackers could target more victims (Section 3) or potentially employ more sophisticated segment prediction to increase the effectiveness of the attack.

We also re-ran the same set of experiments with a set of hubs in place of the switch, effectively removing queuing from the system. The times to download the 100MB file while under attack were reduced to 5 times and 4.5 times the baseline for LAN and WAN attackers, respectively. In other words, having queuing on the bottleneck link significantly *increased* the damage from the attack. We surmise this is because the opt-ack attacker used  $mss = 536$  and the non-malicious client, since it was on local Ethernet, used  $mss = 1448$ . Once the queue was full, the switch could service two of the attack packets before there was room for a legitimate (i.e. destined to the non-malicious client) packet. Effectively, the higher rate of smaller packets caused the switch to drop more non-malicious/legitimate packets. Removing the queue from the system reduced the amount of dropped legitimate packets, therefore increasing non-malicious throughput.

### 5.2.2 Amplification Factors

To evaluate the potential effectiveness of the distributed opt-ack attack, we measure the amount of traffic that our implementation code can induce in a single victim. In this experiment, we use the LAN attacker, as above, to attack various operating systems including GNU/Linux 2.4.24, Solaris 5.8, Mac OS X 10.2.8, and Windows XP with service pack 1. For this experiment, instead of a web server, each victim ran a program that streamed data from memory. This was done to remove any potential application-level bottlenecks from the experiment. As above, the attacker used parameters  $\beta = 16000$ ,  $mss = 536$ , and  $wscale = 0$ . We measured the bandwidth in one second intervals using a custom tool written with the libpcap library. Each experiment in Table 2 was run 10 times, averaged, and is shown as an amplification factor of the attacker’s used local bandwidth.

We believe that the variation in amount of flooding by OS is due to the lack of sophistication of our attack implementation. The amplification factor for Linux is 251.6 times the used bandwidth, which translates to  $251.6/1336$  or approximately 18% of the theoretical maximum traffic,  $\mathcal{T}_{max}$ . This low number is in part because the implementation sends four ACKs per window (as described in Section 3), which alone limits the attack to 25% of  $\mathcal{T}_{max}$ .

## 5.3 Performance of Skipped Segments Solution

In the final experiment, we evaluate the efficiency of our proposed randomly skipped segments solution. Specifically, we measure the time for a non-malicious client on the LAN with selective acknowledgments (SACK) enabled to download a 100MB file from the server with and without the fix and with various threshold val-

Experiment	Time(s)	Deviation	%
Unfixed	89.136	0.007	100%
Fixed: 10-20	89.623	0.980	99.457%
Fixed: 1-200	89.158	0.0234	99.975 %
Fixed: 100-200	89.167	0.0256	99.965 %

**Table 3: Time to Download a 100MB File for Various Fix Options - SACK Enabled**

ues for the fix. The download times were measured with the UNIX *time* utility. Each experiment was run ten times, the results were averaged and are presented in Table 3 with SACK enabled. The two numbers in the first column refer to the range from which the randomly skipped segment was chosen. Results with SACK *disabled* are presented in [22], but are summarized here.

The results show that the performance hit from the proposed fix is negligible for most parameters. Even when we chose the threshold to be intentionally inefficient, i.e., skipping a segment every 10 to 20 ACKs, the fix maintained 99.457% efficiency. We found that varying low end of the range had little effect when combined with SACK, but made a 1% difference with SACK disabled. We believe the performance loss from skipping segments every 100-200 ACKs, i.e., less than 0.1% with or without SACK, is an acceptable price for defeating this attack.

## 6. RELATED WORK

There are two works directly related to opt-ack, which we address below. Due to space considerations, a more complete treatment of related work can be found in [22].

### 6.1 Misbehaving Receivers

As previously mentioned, Savage et al.[21] discovered the opt-ack attack as a method for misbehaving receivers to get better end-to-end performance. While they suggest that opt-ack can be used for denial of service, they did not investigate the magnitude of the amplification the attack can achieve. As a result, their cumulative nonce solution to the opt-ack attack does not consider global deployment as a goal. In this work, through analysis and implementation, we have shown that opt-ack is a serious threat. Further, we have engineered an efficient solution that does not require client-side modification, and thus is more readily deployable.

### 6.2 Reflector Attacks

In [16], Paxson discusses a number of attacks where the initiator can obscure its identity by “reflecting” the attack off non-malicious third parties. As a general solution, Paxson suggests upstream filtering based on the attack signature with the assumption that it is not possible to overwhelm the upstream filter with useless data. The work specifically mentions that if the attacker is able to guess the ISN of the third party, it is possible to mount a blind opt-ack attack against an arbitrary victim. No analysis is made of the amount of the amplification from the opt-ack attack, nor is it immediately clear what filter rules could be applied to arbitrary TCP data.

## 7. DISCUSSION AND CONCLUSION

We have described an analysis of the opt-ack attack on TCP and demonstrated that amplification from the attack makes it dangerous. We have also engineered an efficient skipped segments defense against attacks of this type that allows for incremental deployment. The opt-ack attack succeeds because it violates an un-

derlying assumption made by the designers of TCP: that peers on the network will provide correct feedback. This assumption holds when clients are interested in receiving data, since false feedback will usually lead to worse end-to-end performance. However, the opt-ack attack shows that if malicious nodes do not care about data transfer integrity, they can cause widespread damage to other clients and to the stability of the network.

Since opt-ack violates an underlying assumption upon which TCP is based, we believe a proper solution for the opt-ack attack involves changing the TCP specification. Although new features can be added to TCP (e.g., cumulative nonces) to ensure the receiver TCP is in fact receiving all of the segments, this type of solution is difficult to deploy because it requires client modification. The skipped segment solution presented here requires modification of only high capacity servers, and is thus more readily deployable. In this paper, we have described different mechanisms that can be used to defend against opt-ack attacks. We recommend a specific change to the TCP specification that we have shown to be easy to implement, efficient for fast connections, and which does not burden resource-poor hosts.

**Acknowledgments.** We thank the anonymous referees of the CCS2005 program committee, as well as David Levin, Neil Spring, and Virgil Gligor for their helpful comments and suggestions.

## 8. REFERENCES

- [1] Gnutella Home Page. See <http://gnutella.wego.com>.
- [2] See [www.kazaa.com](http://www.kazaa.com).
- [3] RFC793: Transmissions Control Protocol, September 1981.
- [4] <http://www.akamai.com/en/html/technology/overview.html>.
- [5] T. H. P. R. Alliance. Know your enemy: Tracking botnets. <http://www.honeynet.org/papers/bots/>.
- [6] M. Allman and V. Paxson. On Estimating End-to-End Network Path Properties. pages 263–274. SIGCOMM, 1999.
- [7] S. Floyd and K. Fall. Promoting the Use of End-to-End Congestion Control in the Internet. *IEEE/ACM Transactions on Networking*, 7(4):458–472, 1999.
- [8] V. Jacobson. Congestion Avoidance and Control. *ACM Computer Communication Review; Proceedings of the Sigcomm '88 Symposium in Stanford, CA, August, 1988*, 18, 4:314–329, 1988.
- [9] N. Leibowitz, M. Ripeanu, and A. Wierzbicki. Deconstructing the kazaa network. In *3rd IEEE Workshop on Internet Application*. IEEE, 2003.
- [10] J. Liang, R. Kumar, and K. W. Ross. The KaZaA Overlay: A Measurement Study. <http://cis.poly.edu/~ross/papers/KazaaOverlay.pdf>.
- [11] M. Lichtenberg and J. Curless. DECnet Transport Architecture. *Digital Technical Journal*, 4(1), 1992.
- [12] <http://www.limewire.com/english/content/netsize.shtml>.
- [13] M. Mathis, J. Mahdavi, S. Floyd, and A. Ramanov. RFC2018: TCP Selective Acknowledgment Options, October 1996.
- [14] D. Moore, V. Paxson, S. Savage, C. Shannon, S. Staniford, and N. Weaver. Inside the Slammer Worm. See <http://www.caida.org/outreach/papers/2003/sapphire2/>.
- [15] Nanog email: Dos? <http://www.merit.edu/mail.archives/nanog/2003-01/msg00594.html>.
- [16] V. Paxson. An Analysis of Using Reflectors for Distributed Denial-of-service attacks. *ACM Computer Communications Review (CCR)*, 31(3), July 2001.
- [17] Planet lab. <http://www.planet-lab.org>.
- [18] K. Ramakrishnan, S. Floyd, and D. Black. RFC3168: The Addition of Explicit Congestion Notification (ECN) to IP.
- [19] D. Reed. "Small TCP Packets == very large overhead == DoS?", July 2001. See <http://www.securityfocus.com/archive/1/195457>.
- [20] S. Saroiu, K. Gummadi, and S. Gribble. Measuring and Analyzing the Characteristics of Napster and Gnutella Hosts. volume 9, pages 170–184. Multimedia Systems, Springer-Verlag, 2003.
- [21] S. Savage, N. Cardwell, D. Wetherall, and T. Anderson. TCP Congestion Control with a Misbehaving Receiver. *Computer Communication Review*, 29(5), 1999.
- [22] R. Sherwood, B. Bhattacharjee, and R. Braud. UMD-CS TR #4737 - Misbehaving TCP Receivers Can Cause Internet-Wide Congestion Collapse. Technical report, University of Maryland, Computer Science Department, 2005.
- [23] S. Staniford, V. Paxson, and N. Weaver. How to Own the Internet in Your Spare Time. USENIX Security Symposium, 2002.
- [24] W. R. Stevens. RFC2001: TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms, January 1997.
- [25] P. Watson. Slipping In The Window: TCP Reset Attacks. See [http://www.osvdb.org/reference/SlippingInTheWindow\\_v1.0.doc](http://www.osvdb.org/reference/SlippingInTheWindow_v1.0.doc).