

Opportunities and Limits of Remote Timing Attacks

SCOTT A. CROSBY

Rice University

and

DAN S. WALLACH

Rice University

and

RUDOLF H. RIEDI

Rice University¹

Many algorithms can take a variable amount of time to complete depending on the data being processed. These timing differences can sometimes disclose confidential information. Indeed, researchers have been able to reconstruct an RSA private key purely by querying an SSL web server and timing the results. Our work analyzes the limits of attacks based on accurately measuring network response times and jitter over a local network and across the Internet. We present the design of filters to significantly reduce the effects of jitter, allowing an attacker to measure events with 15-100 μ s accuracy across the Internet, and as good as 100ns over a local network. Notably, security-related algorithms on web servers and other network servers need to be carefully engineered to avoid timing channel leaks at the accuracy demonstrated in this paper.

Categories and Subject Descriptors: C.2.0 [**Computer-Communication Networks**]: General—*Security*; C.2.5 [**Computer-Communication Networks**]: Local and Wide-Area Networks—*Internet*

General Terms: Security, Measurement

Additional Key Words and Phrases: Information leakage, Jitter, Timing attacks

1. INTRODUCTION

Security researchers have studied a number of remote timing attacks, principally against cryptographic algorithms. If an attacker can precisely time cryptographic operations, the attacker may be able to solve for the cryptographic key. There has been significant interest in these attacks. Brumley and Boneh [2004] showed that such attacks were practical,

¹On leave at the School of Engineering, EIA, Fribourg, Switzerland

This research was supported by NSF grants #CNS-0509297 and ANI-0338856 as well as generous support from Microsoft and Schlumberger. The authors wish to thank Alan Cox, Eugene Ng, Scott Rixner and the anonymous referees for their helpful comments.

Authors's address: Dept. of Computer Science, Rice University, 6100 Main St. MS 32, Houston, TX 77005.

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 20YY ACM 0000-0000/20YY/0000-0001 \$5.00

i.e., an attacker could measure the response-time variances of a secure web server with carefully chosen input and, after collecting enough samples, could derive that server's RSA private key. Likewise, it has been shown that an IMAP password could be extracted from a TLS/SSL channel across a network consisting of two switches and a firewall by measuring a 2ms difference in response time [Canvel et al. 2003].

Brumley and Boneh based their attack on roughly 1.4 million queries which they found to be effective on a local area network yet ineffective across the Internet to derive the server's private 1024-bit RSA key. Clearly, additional network hops will increase the latency for packets to travel from one host to another. To an attacker trying to mount a timing attack, latency differences are irrelevant because the attacker is only interested in measuring the *differences* in latency across measured events. However, additional network hops may also add *jitter* (i.e., random noise) to the measured latency. An attacker's goal is to make multiple timing measurements and hopefully smooth out the jitter to recover the actual time difference.

Timing attacks have broad relevance beyond protecting cryptographic keys. Consider algorithmic complexity attacks [Crosby and Wallach 2003], where an attacker tries to induce algorithmic worst-case behavior in a program by sending carefully-chosen inputs that might, for example, cause every insert into a hash table to collide, causing expected $O(1)$ operations to consume their worst-case $O(N)$ running-time. One proposed solution to such attacks is to hide important details about the parameters used by internal algorithms. For example, several systems have replaced a deterministic hash function with a keyed but non-cryptographic hash function (e.g., universal hashing [Carter and Wegman 1979] or Jenkin's hash [Jenkins 1995]). If an attacker can measure a server's response time with enough accuracy to determine if a collision occurred, then the attacker might be able to derive the key.

Timing attacks against some algorithms will require more precision than others. This paper aims to quantify the precision that an attacker might hope to achieve in discriminating between two events, on a remote computer, that take slightly different amounts of time to run. This paper will present the results of extensive measurements both on our local network and across the Internet.

Section 2 describes our attacker model and summarizes our results. Section 3 describes our experimental infrastructure and how we collected measurements. Section 4 describes our network model. Section 5 provides a statistical analysis of our observed network jitter. Section 6 consists of our simulation-driven study of how well an attacker might be able to perform a remote timing attack using statistical hypothesis testing. Section 7 presents related work. We present our conclusions in Section 8.

2. ATTACK MODEL AND RESULTS

We consider a simplified situation where the attacker can transmit two different requests to the target server that either take the same or different time to complete as a function of the server's secret. We assume that knowing whether or not they take the same time will divulge something about the secret. (We assume the attacker knows everything else about the target machine, including its hardware and software configuration. The attacker only lacks knowledge of the internal secret.) The attack then reduces to inferring such a difference in computing time with high reliability. Our model is powerful enough to represent the Brumley and Boneh [2004] OpenSSL attack.

To be more precise, a query is transmitted to the target machine, upon which the target performs some task which requires a *processing time* which the attacker would like to infer. However, the attacker can only measure the *response time*, i.e., the time from when the query is transmitted to when the reply is received.

Our goal is to *identify the smallest difference between two processing times* that can be reliably detected by an attacker given a reasonable number of measurements. Of course, the resolution with which an attacker can observe these differences will be a function of how many samples the attacker takes, how much random perturbation, called *jitter*, is introduced in the response time by the network and how effective the attacker can be at filtering that jitter.

Unfiltered jitter is a measure of resolution based on statistical techniques and is suggestive of the timing difference that is measurable. We also measure resolution empirically by simulating an attack and identifying the minimum distinguishable timing difference. We now summarize our research results with forward references to the subsections where they are discussed.

Does network latency follow a Gaussian distribution? No. The distribution of response time is a highly skewed distribution. (See Section 5.1.)

How is unfiltered jitter measured? Unfiltered jitter estimates the amount of residual noise which limits the resolution of a remote timing attack. Measurements are *filtered* to a single value that is supposed to be correlated with the remote processing time. As we know the actual processing time, we can check the strength of the correlation. Good filters should have high correlation. *Unfiltered jitter* measures the lack of correlation. (See equation 10 in Section 5.4.)

Isn't mean (or median) the best way to filter measurements? If the response times were distributed in a Gaussian fashion, then mean or median would be excellent filters. With the non-Gaussian distributions we see in practice, low-percentile filters tend to significantly outperform the mean or median. (See Sections 5.1 and 5.5 for details.)

Then surely the minimum response time should be an excellent choice. Contrary to expectations, the minimum response time is *not* the least noisy signal. Low percentile filters exhibit significantly less noise than the minimum response times. (See Section 5.5.)

Is there a correlation between unfiltered jitter and network latency or hopcount? We observed no significant correlation between unfiltered jitter and either latency or the number of network hops. (See Section 5.6.)

Does the CPU of a machine affect unfiltered jitter? Yes, we found that our Pentium 4 measurement host introduced artifacts. For processing times between 100ns and 80 μ s, 40% of the measurements had a constant 180 μ s response time. (See Section 5.7.)

Does the network card affect unfiltered jitter? Yes, we found that our Intel Gigabit Ethernet card had 10 to 30 times more unfiltered jitter than our generic onboard Ethernet adapter, depending on whether interrupt coalescing was enabled. (See Section 5.8.)

What is the best statistical discrimination test? Our best test for distinguishing whether two sets of messages had the same response time compared the results of two low-percentile

filters. (See Section 6.4.)

At what empirical resolution can an attacker time a remote host? The resolution an attacker can time a remote host depends on how many measurements they can collect. Our simulated attacker using statistical hypothesis testing was able to reliably distinguish a processing time differences as low as 200ns and 30 μ s with 1000 measurements on the LAN and WAN respectively with. (See Section 6.)

How much extra noise is introduced by application load? With 1000 samples, the addition of application load from Apache introduces only 1 μ s of jitter, much less than the 20 μ s WAN jitter. (See Section 6.5.)

3. EXPERIMENTAL SETUP

We ran tens of millions of timing experiments, both in our lab and over the Internet. Our system implemented a simple UDP ping-pong protocol where, for each measurement, a client sends a message to the server containing the specific amount of time the server should pause before replying. The server waits for the requested amount of time and then responds.

Upon receiving the response, the client logs the processing time requested of the server and the observed response time, then waits a random delay before sending the next request. This delay, averaging 20ms, avoids synchronization artifacts between the client and server. Furthermore, each client performed its trials in a random order. If no response is received within one second, the client assumes that a packet was dropped and repeats the measurement.

In all of our datasets the server machine was dedicated to the task, while we ran clients as background tasks on other machines. Clients would only make one request at a time, allowing them to run without disturbing the machines' users. In effect, we have busy and idle clients querying an idle server instead of mostly-idle clients querying a busy server, as we would expect when a server is under attack. We show measurements of loaded servers and discuss this issue further in Section 6.6.

3.1 Clock calibration

We perform our measurement and delay timings by using the CPU cycle counter. This is the most precise time-source we have available in the system, although non-determinism in CPU behavior reduces the resolution of the clock. When we are simulating processing delays, we read the clock-cycle counter in a busy-loop without a synchronizing instruction. When measuring the response time, we read the cycle counter before having the kernel send the request and right after the kernel returns with the reply message. Beyond this, we rely on our filtering strategy to yield additional measurement accuracy.

Cycle counters count in cycles, which we must convert to nanoseconds, requiring our experimental harness to estimate the clock frequency of each machine. Rather than trying to tightly synchronize the clocks of these machines, perhaps with NTP [Mills 1991; 1992], we perform a one-second calibration of the cycle counter against the system clock to give an approximate solution, within 1% accuracy. Our network model includes a correction for clock skew. Each machine measures time independently, in nanoseconds, and we reconcile the differences in post-processing. We determine the clock skew by the slope of a least-squares linear fit of the delays requested with the delays measured. This process is

Table I. Dataset statistics.

Dataset	#hosts starting	#hosts surviving	total #trials	# samples	Duration
L	8	8	9.0M	27k	1 day
W_A	75	51	112.5M	40k	8 days
W_B	103	37	85.6M	30k	6 days
W_C	136	91	68.4M	13.5k	6 days
W_D	124	3	179.6M	68k	14 days

described further in Section 5.4.

3.2 Collected datasets

We collected five primary datasets. One dataset was collected over a LAN and the other four datasets were collected over the Internet. In each dataset, we measured the same $M = 46$ distinct processing times on the server ranging from 100ns to 8ms. Table I summarizes these measurements, including the number of hosts that were involved at the start of the experiment, the number of hosts that were left when we finished the experiment, the total number of measurements in the dataset, and the maximum number of samples per host per processing time we collected.

Dataset L consists of 8 clients on our LAN. For each processing time, we collected 27,000 samples. Datasets W_A , W_B , W_C , and W_D were collected over the Internet, ranging from 13,000 to 68,000 measures per host per processing time. Dataset W_A was collected in the spring of 2004 and datasets W_B - W_D were collected in the spring of 2005.

To achieve a broad sampling of Internet hosts, we used PlanetLab, an open, globally distributed platform for developing, deploying and accessing planetary-scale network services [Chun et al. 2003]. Unfortunately, PlanetLab hosts are not very reliable and many hosts failed during the experiment. We are unable to restart a failed host because the clock calibration would not match. These host failures were particularly apparent in datasets W_B - W_D .

We would have preferred to use dataset W_A because PlanetLab hosts were much more reliable and less overloaded when it was collected. Unfortunately, a concurrency bug in our earlier data collection system corrupted about .002% of the measurements in dataset W_A and about .0003% in dataset L . We fixed this problem for subsequent measurements, and our post-processing removed all suspicious samples we could identify. Nonetheless, we will present results for each dataset.

4. NETWORK MODEL

It would be impossible to isolate and measure every possible contribution to network latency and jitter. Between an attacker’s machine and the target machine, there may be any number of network bridges, hubs, switches, firewalls, and routers. Each of these may delay packets, drop packets, or suffer internal contention. Furthermore, if packets are arriving faster than they can be forwarded, a router will attempt to queue the packets and send them out later. As the load varies, so will the latency and jitter accumulated by packets as they pass through the network device.

In addition, the end-hosts will introduce their own jitter as a result of application load, virtual memory pressure, and network packet processing. Whereas the attacker may be able to dedicate a computer to the sole purpose of time measurement, and thus reduce the

attacker’s contribution to jitter, the target machine is likely to be running a general-purpose operating system and supporting a non-trivial workload.

We use an abstract model of the server and the network. The server is assumed to run at least two different *tasks* which have different processing times and all requests for the same task have the same processing time. We model the latency of a round-trip communication channel between one pair of hosts as

$$\text{responseTime} = a \cdot \text{processingTime} + \text{propagationTime} + \text{jitter} \quad (1)$$

with the following five assumptions:

- (1) *responseTime* is the measured round trip time on the network.
- (2) *a* accommodates clock skew and is constant over all requests for all tasks. This is estimated in our analysis independently for each host.
- (3) *processingTime* is *constant* for all requests for the *same* task. In our dataset, this corresponds to the time the remote hosts delays before sending a reply to a ping-pong.
- (4) *propagationTime* is constant over all requests, for *all* tasks. This is the ‘average latency’ and is estimated in our analysis independently for each host.
- (5) *jitter* is the jitter term and is identically and independently *distributed* for all requests and tasks.

Thus, the jitter term absorbs all randomness introduced by network conditions, load on the target machine and any other source.

More formally, the M different *known* processing times are denoted by $t[1], \dots, t[M]$. The true, but unknown statistical distribution of response times to a query m with a processing time of $t[m]$ is denoted by $R[m]$. Thus, for each host pair and for each query m , $r_1[m] \dots r_N[m]$ denote the N collected measurements which may also be thought of as samples taken from $R[m]$. Our channel model (1) can be restated more formally as:

$$\forall_{m \in \{1 \dots M\}} \quad r_n[m] = a \cdot t[m] + b + \epsilon_n[m] \quad (2)$$

where b is the propagation time, a is the correction for clock skew and $\epsilon_n[m]$ is the jitter in the measurement $r_n[m]$.

By our channel model, the random variables $\epsilon_n[m]$ are all of the same distribution and are statistically independent (for all n and m). In other words, the distributions of $R[m]$ for different m are identical up to a shift term and the $r_n[m]$ constitute independent samples of these distributions.

Our model setup implicitly assumes that the application load is stationary and that the network will introduce the same random perturbations for all requests. This is reasonable if all response messages have the same payload and if routing is stable over the time of the measurements. An attacker may be limited to the number of measurements they can collect due to route changes. Combining measurements across route changes would require detecting when those changes occur and determining the propagation time difference between each set of samples. (Our own work assumes that the network routes are constant across the duration of our experiments.)

We repeat our analysis for each host pair in each dataset, treating each of them as a separate channel with unique jitter, propagation time, clock skew, and response time distribution.

The definition of jitter The splitting of the response latency into the sum of propagation and jitter is for convenience and does not affect methodology or performance of

the attack. Indeed, since only *differences* in response times need to be inferred, and not actual estimates of the response times themselves, introducing the constant shift term $b = \text{propagationTime}$ is irrelevant. In particular, we do not require an accurate estimate of the propagation time as long as we use the same constant value throughout the attack. Although strictly speaking, jitter is always an additional positive latency, we consider propagationTime as an estimate of the propagation time plus the *average jitter*, allowing *jitter* to represent positive and negative deviations from that average.

The use of clock skew An attacker need not be concerned with clock skew on a remote host. They only need to reduce their own clock skew enough to avoid gross errors in their measurements. Ignoring this skew creates a relative error of $(a - 1) \approx 1\%$ for an attacker, such as the one we simulate in Section 6, which is of no importance.

In contrast, an accurate assessment of clock skew between our client hosts and our server host is critical in our meta-analysis of Internet jitter. Our analysis of jitter depends upon tighter clock calibration than would be necessary to perform a timing attack. If the server and client had their clock skew miscalibrated by 1%, this would cause a systematic bias between the processing time we actually measured and the processing time we believed we measured. For example, if we use the client’s clock as a baseline, with the server’s clock slightly miscalibrated to run 1% faster, the processing time delay performed by the server, as measured by the client’s clock, would be too low, and the reply message would be sent too early.

Allowing the least-squares fit, used to estimate unfiltered jitter, to also estimate this clock skew a compensates for this bias and prevents us from radically overestimating the jitter. For example, where our current system estimates a 100ns jitter, the lack of clock skew compensation would have yielded an $80\mu\text{s}$ estimate of jitter — three orders of magnitude higher than the real jitter.

5. STATISTICS OF THE RESPONSE TIME DISTRIBUTION

We first look at the response time distribution. We examine techniques to filter it and estimate the unfilterable jitter.

5.1 Response time distribution

The distribution of jitter is not Gaussian. It is highly skewed distribution with two modes. In Figure 1 we plot a histogram of the probability density function (PDF) of the response times for one long-uptime host chosen from dataset W_D . The response time is clearly asymmetric and non-Gaussian and includes two obvious modes and an exponentially distributed tail. In Figure 2 we plot the CDF corresponding to this host. The steep slope at about 52ms is the least varying part of the distribution and occurs in the 5th-15th percentile response time. Most other hosts have a similar response time distribution, with a few unusually fast responses, and many have more than one mode.

5.2 Measurement quantiles

We define the i -th percentile (or quantile) q_i of a distribution R to be the smallest real number $q_i[m]$ such that $P[R \leq q_i] = i/100$. The 50-th percentile is the familiar *median* and the 0-th percentile is the minimum response time.

The true quantile q_i is unknown because it depends on the true distribution R . We can compute an *estimator* \hat{q}_i , which is an empirical estimate of q_i derived from our measure-

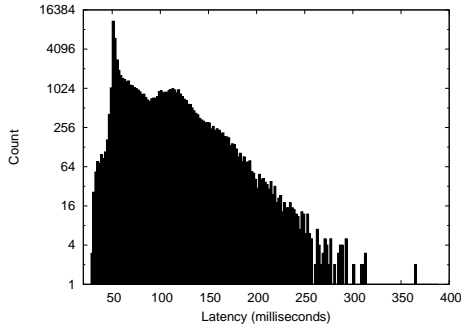


Fig. 1. Histogram of response times for 66k measurements of a 1ms processing time from a host from dataset W_D .

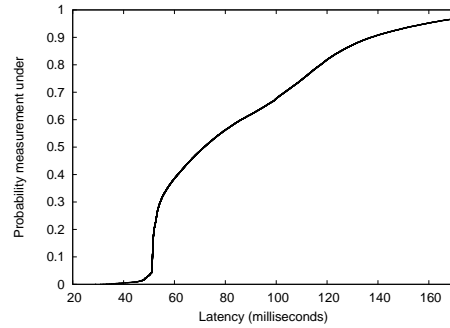


Fig. 2. CDF of response times for 66k measurements of a 1ms processing time from a host from dataset W_D .

ments by first ranking a set of measurements into

$$r_{(1)} \leq r_{(2)} \leq \dots \leq r_{(N)} \quad (3)$$

and setting $\hat{q}_i = r_{(\lfloor iN/100 \rfloor)}$. The empirical quantiles \hat{q}_i are well known to be weakly consistent under mild conditions, meaning that they converge to the true quantile q_i as the sample size increases.

Rather than show a set of CDFs or histograms for each processing time for each host, which can be seen in Figures 1 or 2, we use quantiles as a graphically powerful alternative for summarizing the response time distribution for a single host. We summarize the response time distribution with *percentile contours*. The empirical percentiles $\hat{q}_i[m]$ are plotted as a function of the sampled server processing times, $t[1], \dots, t[M]$, which are indicated with dots on the x -axis. These percentile contour summaries of the response time distribution are easier to compare than their full histograms.

Figures 3 and 4 summarize the estimated percentiles of the response time distributions we observed for a host in dataset L and W_D respectively. Figure 3 shows that that despite intuition to the contrary, the minimum response time seems to be poorly correlated to the processing time. For the host in Figure 4, the minimum response time is about three times noisier than the 1st percentile. (In Section 5.5 we will elaborate on which quantiles make the best filter over all hosts.)

We can present all $M = 46$ response times on one plot on a log-scale if we first estimate and subtract off the linear fit. (This process is further described in Section 5.4.) Figure 5 shows the response time distribution for large and small quantiles for a host in dataset W_D . The variation of the median and minimum response times across different processing times is visually more than the variation seen in measurements at small quantiles. The mean (not plotted) is similarly noisy. Because the median and mean of our measurements are extremely noisy, parametric inference techniques based on these very classical statistics are unlikely to work well.

Variability of the Empirical Percentile To study the error of this percentile estimation we need to consider the distribution of the estimator itself, i.e., the variability in \hat{q}_i when the estimation is repeated. Fixing a percentile i and setting $k = k(i) = \lfloor iN/100 \rfloor$ we may

write

$$\hat{q}_i[m] = r_{(k)} = a \cdot t[m] + b + \varepsilon_{(k)}[m]. \quad (4)$$

Here, we introduced the random variable $\varepsilon_{(k)}[m]$ which is the k -th ranking sample out of N independent and identically distributed samples from $\varepsilon[m]$. Since under our channel model all $\varepsilon[m]$ are identically distributed for different m , $\varepsilon_{(k)}[m]$ are identical in distribution as well and we conclude that the distribution of $\hat{q}_i[m]$ and $\hat{q}_i[m']$ are identical up to the shift term $a(t[m] - t[m'])$:

$$\hat{q}_i[m] \stackrel{\text{distr.}}{\equiv} \hat{q}_i[m'] + a(t[m] - t[m']). \quad (5)$$

Precision of Percentile Estimation It follows immediately from (4) that the estimation errors $\hat{q}_i[m] - q_i[m]$ themselves have identical distributions, independently of m :

$$\hat{q}_i[m] - q_i[m] = (a \cdot t[m] + b + \varepsilon_{(k)}[m]) - (a \cdot t[m] + c) \quad (6)$$

$$= \varepsilon_{(k)}[m] + b - c \stackrel{\text{distr.}}{\equiv} \varepsilon_{(k)}[m'] + b - c. \quad (7)$$

Two remarks are in order.

First, to study the variability of a percentile estimator \hat{q}_i and the estimation error (6), we need a set of empirically obtained values. Typically, this is done by repeating the same estimation procedure several times. Here, we may exploit that the errors $\hat{q}_i[m] - q_i[m]$ obtained for the *same* percentile but for different values of the response time, i.e., for $m = 1 \dots M$ are all of the same distribution (see (7)). In other words, we may consider these error values as samples of the same distribution $\varepsilon_{(k)} + b - c + a \cdot (t[m] - t[m^*])$. In this context we recall that b is the (unknown) propagation time, a is the clock skew, and that c depends only on i and some arbitrary baseline response time $t[m^*]$, but not m .

Second, since we are interested in *differences* of processing times, and thus in (additive) differences of the response times we observed, the additive constant $(b - c)$ in (7) will cancel out in our inference schemes and be of no importance.

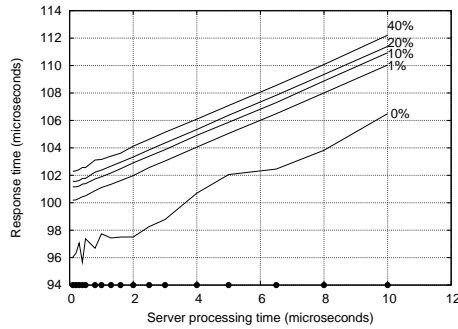


Fig. 3. Response time percentiles as a function of processing time for a local host in dataset L .

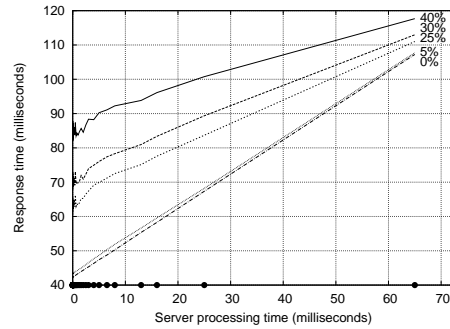


Fig. 4. Response time percentiles as a function of processing time for a remote host in dataset W_D .

5.2.1 *Filtering our data.* Because of the noise in the jitter distribution, particularly at and above the median, we apply a “filter” Γ which reduces the set of measurements $r_1[m] \dots r_N[m]$ to a single number that is hopefully closely correlated with the processing time $t[m]$. A filter, foremost, is a function of the measured response times. Filter design choices are driven mainly by the objective to minimize the variance as a measure of error which in turn impacts the reliability of a decision procedure based on the filter. In general, we may think of this procedure as some sort of de-noising.

While such a filter can not remove all noise, it will in help us understand the unfilterable jitter in our measurements and estimate the resolution an attacker may be able to measure. There are many filters we may choose. In Section 5.4 we describe how we evaluate filters and identify the ‘best’ such filter.

A first simple example of a filter would be the i -th empirical percentile:

$$\Gamma_{q_i}(r_1[m], \dots, r_N[m]) = \hat{q}_i[m] \quad (8)$$

Recall that the filters $\Gamma_{q_{50}}, \Gamma_{q_0}, \Gamma_{q_{100}}$ pick the *median*, *minimum* and *maximum* response times, respectively.

We tried three types of filtering strategies beyond simple percentile filtering but they performed no better than the best percentile filter. The best filter of each type was identified through brute force: We ran several hundred and recorded the best instance of each type. The four types of filtering we applied are:

- (1) **Percentile** This filter returns the x th quantile or percentile. We tried a hundred different percentile values, varying from 0% to 70%, with most of the percentile values less than 10%.
- (2) **Peak** This filter identifies modes in our measurements. Samples are first sorted, and then a window of a fixed width is moved across the sorted list. When the difference between the maximum and minimum sample within the window is minimized, that position of the window contains the highest density of samples. The median measurement within the window is reported as the result of the filter. We considered 25 different window widths, ranging from 1% to 70% of the collected samples. Even in the best case, our percentile filter outperformed the peak filter.
- (3) **Average Range** Inspired by the filter used in Brumley and Boneh’s SSL attack, we computed a histogram over the samples and computed the average of the samples between the i -th and j -th percentiles. We tested this filter over 96 different percentile ranges. We used percentiles ranging from 0% to 40%; we mostly examined ranges below the 5th percentile. In general, the best instances of this filter performed within a few percent of the best percentile filter.
- (4) **Percentile Smoothing** This is an attempt to improve our estimate \hat{q}_i of the real percentile q_i . Instead of computing the estimate from our measurements directly, we divide the measurements into k disjoint subsets, compute a separate estimate for each subset, and then average those estimates. We tested this filter for $k \in \{4, 10, 20\}$ and for 21 different percentiles. This filter performed slightly worse than our best percentile filter.

5.3 Filtered measurements

Figures 3 and 4 show several percentile filters. For dataset L , virtually all of the percentiles we plot show a smooth relationship to processing time, except for $\Gamma = q_0$, the minimum re-

sponse time. For dataset W_D , it is visually clear that the 25th percentile is much noisier than the 5th percentile. We note, for the particular client of Figure 3, and for plotted processing times ranging from 100ns and 8ms, 39 percent of the 27,000 samples for each processing time occurred within a $2\mu\text{s}$ window, while the minimum response time was about $4\mu\text{s}$ faster. The minimum response time is clearly much noisier than the other percentiles.

In Figure 3 we also observe that indeed the first and tenth empirical percentiles follow a linear dependence on processing time with high accuracy. This is suggestive that the first and tenth percentile are better filters than higher percentiles.

In Figure 4, we show a similar plot for a typical host from dataset W_D . First, note that the scale for this plot is in milliseconds and note that the 25th percentile looks noisier than the 40th percentile of dataset L . This indicates more variability in the jitter in PlanetLab measurements than in the LAN measurements. For fine-grained measurements of response time, up to 75% or more of Internet measurements may be unacceptably noisy. This plot also appears to show non-parallel percentile contours. This is only an artifact caused by the non-uniform spacing of processing times where the 25ms and 65ms measurements are connected with a straight line.

5.4 Verifying channel model

We summarize the jitter distribution as a single filtered value and verify our channel model by how close the filtered measurements match the expected linear relationship. Let $\Gamma(r_1[m] \dots r_N[m])$ be a filter on the measurements $r_1[m] \dots r_N[m]$ for all $m \in 1 \dots M$, reducing them to a single value for each of the processing times $t[1] \dots t[M]$. We verify our channel model by using a least squares linear fit between the M different processing times and the M different filtered response times. Under the network model, the following relation should hold for $m = 1 \dots M$:

$$\Gamma[m] = \Gamma(r_1[m] \dots r_N[m]) = a \cdot t[m] + b \quad (9)$$

The least square fit will provide estimations of the clock skew a , the propagation time b and the variance σ^2 of the estimator, giving us a *measure of confidence* that the data follows a linear relationship. To this end, we compute the unique values \hat{a} and \hat{b} which minimize the average square deviation s_i^2 of $\Gamma(r_1[m] \dots r_N[m])$ from the linear relation. In other words:

$$s^2 = s_i^2 = \frac{1}{m} \sum_m (\delta[m])^2 \quad \text{where} \quad \delta[m] = \Gamma(r_1[m] \dots r_N[m]) - \hat{a} \cdot t[m] - \hat{b}. \quad (10)$$

Explicit formulas for \hat{a} and \hat{b} are quickly derived using standard calculus.¹

For our best filters, s^2 is very small compared to the processing time differences which leads us to accept the linear relation between processing times and percentiles up to a small random error and *validates* our channel model. Furthermore, $\sqrt{s^2}$ which we call *unfiltered jitter* measures the effectiveness of a filter. The filter with the smallest unfiltered jitter is the most effective filter — the one with the best linear fit to the channel model.

¹Setting $\bar{t} = (1/M) \sum_m t[m]$, $\bar{q} = (1/M) \sum_m \Gamma[m]$, $\bar{tq} = (1/M) \sum_m t[m] \cdot \Gamma[m]$ and $\bar{t}^2 = (1/M) \sum_m t^2[m]$ we have

$$\hat{a} = \frac{\bar{tq} - \bar{t} \cdot \bar{q}}{\bar{t}^2 - (\bar{t})^2} \quad \text{and} \quad \hat{b} = \bar{q} - \hat{a} \bar{t}$$

Table II. Mean and median, among all hosts within each dataset, of the unfiltered jitter.

Dataset	W_A	W_B	W_C	W_D	L	L'
Mean unfiltered jitter (μs)	11.7	14.1	19.0	20.9	2.1	.218
Median unfiltered jitter (μs)	6.5	7.2	6.7	7.4	.050	.046

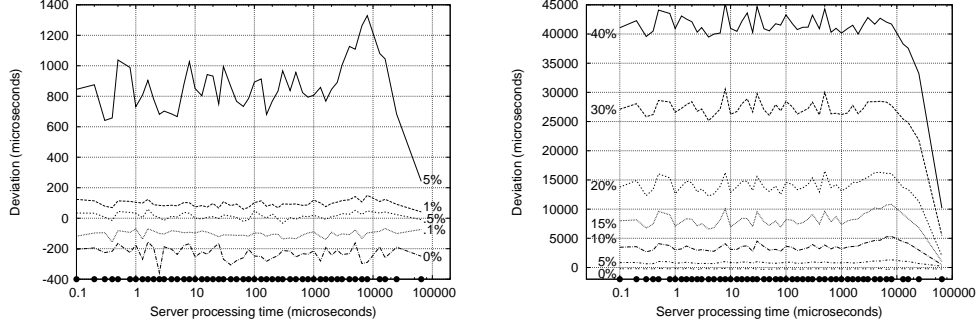


Fig. 5. Difference from ideal estimator as a function of processing time for the same host in dataset W_D , showing small quantiles in the left plot and large quantiles in the right plot. Dots on the X-axis denote measured processing times.

Unfiltered jitter is also suggestive of the resolution that an attacker can distinguish. We may assume that our *filtered* measurements are Gaussian distributed since the filter amounts to an averaging of independent original observations or extracting a percentile. We can use the rule of thumb that a Gaussian distributed sample has less than a 5% probability of being more than 1.7 standard deviations greater than its mean. Thus a timing difference must be at least $3.4 \cdot s$ to be distinguishable with a 5% false positive and 5% false negative rate. In Section 6 we empirically measure an attacker’s resolving power based on our network trace data.

In addition to this validation of our channel model, the least square fit (10) provides *simultaneous* estimates of the clock skew via \hat{a} , and propagation time via \hat{b} . An estimated value $\hat{a} \approx 1$ provides further confirmation of the channel model (1).

Some care has to be taken in the interpretation of \hat{b} . The least squares fit minimizes the sum of the squares of $\delta[m]$, some of which will be positive and some negative. Therefore, we define *define* b to be the *average propagation*, rather than the minimal propagation; jitter may be negative and then the average of $\epsilon_{(k)}$ is zero by definition and $\hat{b} \approx b$.

We summarize the unfiltered jitter across each dataset in Table II by giving the average and median unfiltered jitter. Across all of the datasets, the mean unfiltered jitter is noticeably larger than the median unfiltered jitter, indicating that although many hosts have low unfiltered jitter, some are much worse. For example, we created dataset L' based on dataset L , with one outlier removed. The outlier is a host on the other side of a “traffic shaper” which introduced a significant amount of jitter.

In summary, we accept a small value of s^2 as a validation of the channel model and accept the estimated clock skew \hat{a} to be the true skew a . In statistical terms, this procedure is called a model fitting. We do not overfit as we estimate at most two parameters from $M = 46$ different response times.

We also determine if the estimator error remains the same across all processing times. In Figure 5, we plot the difference from the best-matching percentile filter and the actual mea-

surement percentiles for various processing times on a logarithmic scale. We refer to these plots as *deviation contours* as they emphasize the residual noise in that could not be filtered. Observe the lack of a trend for small percentile values, where the absolute estimator error remains about the same across all processing times. This figure also presents larger quantiles from the same host showing significantly more random variation with short processing times but *less* variation in the measurements for the longer processing times. The jitter distribution is *not* the same across all processing times for this host.

Some hosts follow the channel model and have the same jitter distribution across all processing times, as we would expect. Other hosts violate the channel model, including those using an Intel Gigabit Ethernet NIC or an Intel Pentium 4 CPU (See Sections 5.7 and 5.8). For those, our analysis still finds the best filter, although the resolution of the filtering process may be limited.

5.5 Effectiveness of filtering among quantiles

Intuition says that minimum response time should be the ideal filter because network devices can only introduce additional variable latency. Therefore, the minimum response time should have the least noise introduced. In our experiments to identify the filter with the lowest unfiltered jitter we found filters significantly more effective than using the minimum response time, contradicting this intuition.

To better understand this, we examined the relationship between percentiles and the resultant unfiltered jitter. In Figure 6 we plot the unfiltered jitter for 150 different percentile filters ranging from .001th percentile to 70th percentile, averaged across all hosts in their respective datasets. As expected, percentiles over 10% were very noisy. *Every* dataset shows a trend indicating a higher noise as the percentile declines from the 1st percentile toward the 0th percentile or minimum response time. These curves show several local minima and demonstrate that using the minimum response time leads to several times the error of using the empirically best filter.

To examine this further, we performed separate experiments using a laptop and desktop computer connected with either a crossover cable or a network switch. All hosts and the switch were idle during the experiment. Figure 7 plots the unfiltered jitter as the percentile changes for dataset L as well as the dataset using the switch and crossover cable. We note all three curves have at least two local minima that are not at the 0th percentile. For a switch and a crossover cable, the lowest unfiltered jitter is the 0.9th percentile and 12th percentile, respectively. These results confirm that minimum response time is not an accurate filter and that the ideal percentile filter can be difficult to predict a priori.

5.6 Unfiltered jitter versus network distance

In Figure 8, we show scatter plots of network distance, as measured in the propagation delay of the best linear fit, versus the lowest unfiltered jitter for that host, for each PlanetLab host in our experiment. Each dot on these plots corresponds to the results of doing a least-squares linear fit on the filtered measurements for one host. This plot shows how network round trip times correlate with measurement accuracy.

The largest cluster of hosts in the bottom left reflects the large number of PlanetLab hosts within the United States. Interestingly, the international hosts, with significantly longer network latencies, do not have noticeably higher unfiltered jitter. The host with the least unfiltered jitter in dataset W_A , $2\mu s$, was physically located on another continent. As such, physical distance does not necessarily imply much about unfiltered jitter. The hosts

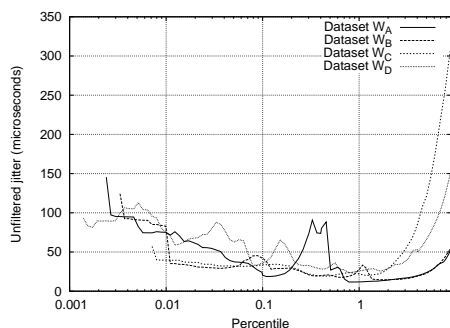


Fig. 6. Unfiltered jitter as a function of filter percentile over PlanetLab.

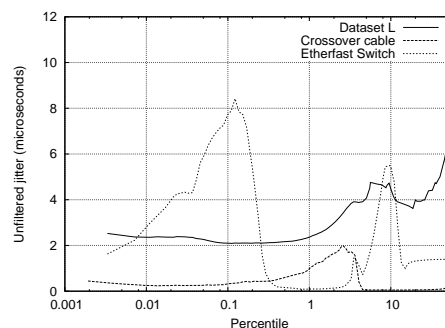


Fig. 7. Unfiltered jitter as a function of filter percentile over the LAN.

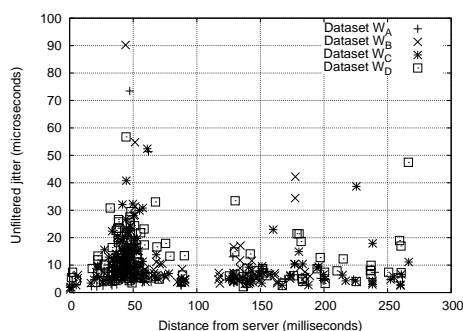


Fig. 8. Scatter plot of unfiltered jitter versus network distance for each host in our PlanetLab datasets.

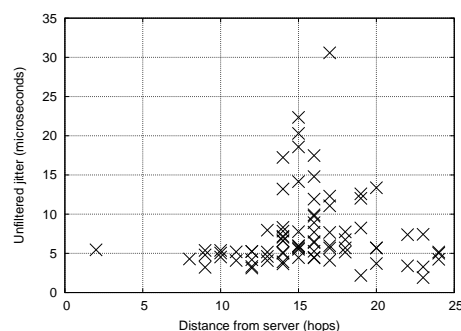


Fig. 9. Scatter plot of unfiltered jitter versus network hop count for 80 hosts in dataset W_D that replied to at least 10k measurements.

with higher unfiltered jitter, notably the outliers in the upper left, may reflect the lack of data we were able to collect from some PlanetLab hosts before they failed (see Section 3.2).

We omit similar figures for dataset L , our LAN measurements. Most hosts were only connected by switches with no router hops and had 50ns of unfiltered jitter. We had one outlier host with $3.5\mu\text{s}$ of unfiltered jitter, which was behind a traffic shaping box.

For dataset W_D , we additionally measured the number of hops from the local system to each PlanetLab host used in the experiment. This allows us to determine whether the number of hops, rather than the latency, influences the unfiltered jitter. Figure 9 plots the unfiltered jitter versus the network hop count for hosts that had at least 10k measurements and a successful traceroute. One PlanetLab host located 2 hops away on the LAN had twice the unfiltered jitter of the best host, located 24 hops away. This plot also shows no correlation between hopcount and unfiltered jitter.

5.7 CPU dependencies

During the course of our experiments we observed unusual behavior when the delay being measured was less than $100\mu\text{s}$. In that configuration, we happened to use a Intel Pentium 4 laptop as the measuring host. Figure 10 demonstrates this problem. 35% of the samples

with a processing time between 100ns and 80 μ s take a constant time to return, regardless of the actual processing time. We initially suspected this was an artifact of a router or switch. In fact, the culprit was the CPU itself.

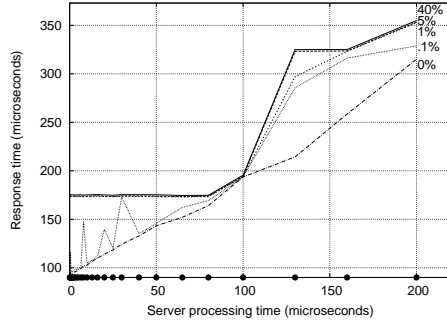


Fig. 10. Response time percentiles as a function of processing time for a Intel Pentium 4 laptop measurement host and AMD Athlon target on a crossover cable.

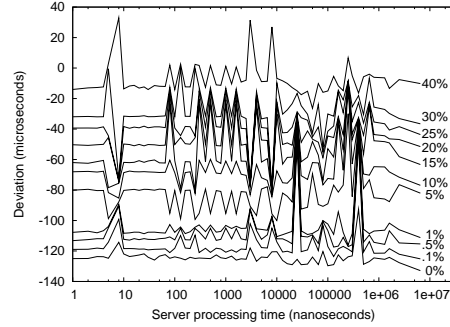


Fig. 11. Deviation contours from ideal straight line for dataset L_A , using an Intel gigabit NIC.

To confirm that, indeed, the Intel Pentium 4 CPU was the source of the problem, we performed 12 timing measurements between four different computers (two models of Intel Pentium 3, a Intel Pentium 4 desktop CPU and an AMD Athlon 1GHz). In these measurements, the artifact occurred only when the Pentium 4 desktop CPU was used in the client host and occurred nowhere else.

Our original Intel Pentium 4 laptop (a 1.8GHz Pentium 4-M) and the new Intel Pentium 4 desktop (a 3.06GHz Pentium 4) have entirely different motherboards, Ethernet devices, and so forth. We conjecture that the artifact may be a consequence of the Pentium 4's power management, putting the computer to sleep at inopportune moments. While a more detailed study of this effect is beyond the scope of this paper, prospective attackers will certainly select and profile their timing host to ensure these artifacts do not occur.

5.8 Network card dependencies

We performed a further experiment to characterize how the choice of networking card (or *NIC*) might impact unfiltered jitter. We ran a new experiment between two machines that each had two networking cards. Both machines had a 2.1GHz AMD Athlon CPU and an Intel 82540EM gigabit ethernet controller. Machine A acted as the sender and has an onboard Via VT6102 100baseT ethernet controller. Machine B acted as the receiver and has an onboard 3c905C-TX/TX-M 100baseT ethernet controller. We connected these machines with either a crossover cable or a generic SpeedStream SS2108 10/100 baseT switch.

Table III summarizes our results on four new datasets where we collected up to 150k measurements for 60 processing times ranging from 1ns to 6.5ms on a 100baseT ethernet network. We found that the higher performance gigabit card, with interrupt coalescing enabled, had 30 times the unfiltered jitter than the generic onboard ethernet card.

There is a caveat to these results. Datasets L_A and L_B violate our network model because their jitter distribution is very different across different processing times. Figure 11 shows

Table III. Unfiltered jitter differences between networking cards.

Experiment	Sender NIC	Receiver NIC	IRQ coalescing	Connection	Measurements	Unfiltered jitter
L_A	Intel	Intel	default	switch	150k	1419ns
L_B	Intel	Intel	disabled	switch	100k	676ns
L_C	Onboard	Onboard	n/a	crossover	100k	56ns
L_D	Onboard	Onboard	n/a	switch	100k	49ns

the deviation contours for the host in dataset L_A and demonstrates the variability in the jitter distribution. These results indicate that some modern high performance servers, with high performance NICs, may be safer from remote timing attacks as a direct consequence of their performance features.

6. SIMULATING ATTACKS

A statistical analysis of unfiltered jitter is useful in that it identifies the lower bound of noise in the system and indicates what resolution is theoretically distinguishable. In this section, we simulate an attacker performing a timing attack and evaluate the empirical resolution that a real attacker could identify. In statistical terminology this amounts to performing a *hypothesis test*.

In hypothesis testing, the statistician tries to show that their observations are statistically significant enough to exclude the default or *null hypothesis*. As our measurement dataset contains millions of measurements of ‘ground truth’, where we know the actual processing time, we may evaluate the empirical effectiveness of different hypothesis testing approaches an attacker may attempt. Our choices of which approaches we try are guided by the results in the last section, where we determined that the mean, median, and measurements above the 10-20th percentile are extremely noisy.

We show and test a framework of statistical hypothesis testing to determine how fine a time difference Δ might be detectable within a reasonable number of measurements, both for LAN and Internet attackers. We find that, with as few as 2000 measurements, it is possible for an attacker to distinguish a 100ns processing-time difference on a LAN with false negative and false positive rates under 5%.

6.1 Classic hypothesis testing

Many hypothesis testing approaches that distinguish whether two measurements sets U and V are from the same distribution or different distribution, such as the Student’s t -test, use a *test statistic* ϕ , computed from the summary statistics of U and V such as the mean, standard deviation, and the sample count. If ϕ exceeds a statistically significant threshold, then the null hypothesis, that U and V are from the same distribution, is rejected.

Of course, any such analysis of hypothesis testing with a single threshold incurs a chance of two types of errors: false positives, when we improperly reject the null hypothesis, and false negatives, when we fail to reject the null hypothesis when we should.

6.2 Empirical hypothesis testing

In our scenario, we are not limited to computing the theoretical effectiveness of a hypothesis test because our datasets contain ground truth. We may simulate an attacker performing empirical hypothesis tests and directly rate the attacker’s effectiveness. Rather than collect a new dataset for hypothesis testing, we reuse our previously-collected raw measurements,

as analyzed in the previous section.

Our simulated attacks follow this procedure: the hypothesis test H is given two sets of N samples of response times $X = r_1[i], \dots, r_N[i]$ and $Y = r_1[j], \dots, r_N[j]$ corresponding to processing times $t[i]$ and $t[j]$ which may or may not differ by Δ . For any given host pair and Δ we wish to discriminate, we can randomly choose data from our network measurements $r_n[m]$ to populate X, Y . We let the null hypothesis be that $\Delta = 0$, i.e., that X, Y are from the same distribution. We then perform the hypothesis test $H(X, Y)$. To *accept the null* means that H believes that X, Y are from the same distribution and to *reject the null* means that H believes it has sufficient evidence that the X, Y are in fact different distributions.

We use the same testing procedure for each hypothesis test to compute the empirical false positive (FP) and false negative (FN) rate for each host pair and Δ . To compute the false negative rate, we perform 200 trials where we populate the queries X, Y from different distributions ($t[j] = t[i] + \Delta$) and count how many times the hypothesis test mistakenly accepts the null. To compute the false positive rate, we perform 200 trials populating the queries X, Y from the same distribution ($t[i] = t[j]$) and count how many times the test mistakenly rejects the null. We summarize these measurements by computing, for each host in question, the smallest Δ with a FN and FP rate below 5% and denote this as *empirical resolution*. If the hypothesis test H_p includes parameters p , we find the ideal parameters by brute force. We try H with all choices and keep the best-performing instance. This simulates a “best-case” attacker that has an oracle that helps it choose optimal parameters.

6.3 Hypothesis testing approaches

We considered four different statistical hypothesis testing approaches:

- (1) **Students t -test or other parametric approaches.** A parametric statistical approach such as the t -test assumes the distribution of the jitter follows a parameterized distribution (e.g. Gaussian, exponential, etc.) and we need only compute parameters of the fit such as the mean and variance. Doing so allows us to design stronger estimation procedures which exploit the assumed structure. However, it also requires us to perform a goodness-of-fit test, to see whether the model is appropriate.
We considered attempting the t -test on our raw measurements. We rejected it and other parametric approaches because our sample mean, median, variance and upper percentiles have high variability, precluding most parametric models as they would be rejected in a goodness-of-fit test. An inference made on the basis of a poor model would be unlikely to be accurate.
- (2) **The Wilcoxon Rank-sum.** The Wilcoxon rank-sum test (also called the Mann-Whitney rank sum test) is a standard non-parametric test to determine if two sets of observations come from the same distribution. It works well, even if the distributions are non-Gaussian, so long as they represented shifted copies of each other. Following our hypothesis test procedure to determine the false positive and false negative rates, we found this methodology significantly under-performed our “box test.”
- (3) **Modified Students t -test.** Our raw measurements have a very skewed, non-Gaussian distribution. We model an attacker who runs the t -test upon filtered measurements instead of the raw measurements.

The modified t -test we apply is parameterized by a filter Γ and a threshold T . We test each of 19 well-performing filters and 12 thresholds. As before, we start with two sets of N samples of response times X, Y . For each of 200 random subsets of size $N/10$

from the N measurements in X, Y , we apply the filter Γ . We estimate the variance and means of $\Gamma(X), \Gamma(Y)$ from these 200 sampled random subsets and apply the standard t -test with a threshold of T on the t -statistic.

Following our testing procedure for hypothesis tests to determine the false positive and false negative rates, we found that this test had half of the resolution of our “box test.”

- (4) **The “Box Test.”** This is our best-performing test. This test exploits our observation that smaller quantiles in our measurements have less noise. Given a measurement set U and two quantiles i, j , we define an interval $[\hat{q}_i(U), \hat{q}_j(U)]$. The test rejects the null if the intervals induced by those quantiles are non-overlapping and in the proper order. Details on this test appear below.

6.4 The box test

Our “box test” is designed to exploit our observation that small percentiles have the least noise. We also argue for this filter design based on its simplicity². The box test is parameterized by two quantiles i, j and, as before, we start with two random subsets X, Y of N measurements.

To perform this test, X, Y , are sorted and two intervals are formed: $[\hat{q}_i(X), \hat{q}_j(X)]$ and $[\hat{q}_i(Y), \hat{q}_j(Y)]$. The test accepts if the intervals $[\hat{q}_i(X), \hat{q}_j(X)]$ $[\hat{q}_i(Y), \hat{q}_j(Y)]$ do not overlap and $[\hat{q}_i(X), \hat{q}_j(X)]$ is before $[\hat{q}_i(Y), \hat{q}_j(Y)]$. We calculated the false positive and false negative rates by following our testing procedure for hypothesis tests.

Since we cannot a priori know which values for the parameters i, j create the most effective statistical test, we perform this experiment for 12,000 i, j pairs, keeping whichever pair i, j has the lowest FN rate while also having a FP rate below 5%. From our measurements, the lower quantiles are the statistically most reliable part of the distribution. Our exhaustive search for optimum parameters chose $i, j < 6\%$ for over half of the hosts. Although a real attacker would not have such an oracle against an unknown target, an attacker could always perform measurements similar to ours against a known machine, hopefully near the target. Also, we prefer to state an upper bound on the capabilities of an attacker.

Figures 12 and 13 demonstrate the typical behavior we observe, where the false negative rate suddenly increases when the time difference is too small to be reliably distinguished. We illustrate with two hosts selected from datasets L and four hosts selected from W_D . Data for other hosts and other datasets reflects a similar pattern: the FN rate remains high until the time difference increases past a critical threshold, at which time the FN rate drops quickly to zero and the discriminator is exceptionally accurate.

Rather than showing similar scatter plots containing every host that we measured, we summarize our measurements by computing, for each host in question, the smallest Δ which we can discriminate with a FN rate below 5%. We call this the *empirical resolution*. We find that LAN hosts can accurately resolve a Δ in the hundreds of nanoseconds, and that PlanetLab hosts can resolve Δ 's around $100\mu s$, with the best hosts resolving $\Delta \approx 30\mu s$.

So far, we have modeled the empirical resolution of an attacker performing 500 measurements. To see how additional measurements would increase the attacker's abilities, we varied the number of simulated measurements from 10 to 10,000. We summarize the histogram over each of the 124 hosts in dataset W_D by its 5 quartiles (0th, 25th, 50th, 75th

²“Complicated computations do not guarantee a valid statistical analysis. Always start ... with a careful examination of the data.” [Moore and McCabe 2004]

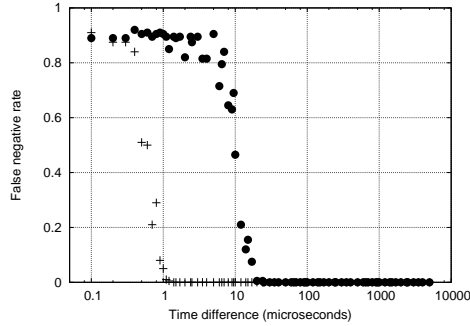


Fig. 12. Scatter plot of the time difference versus the FN rate for two dataset L hosts. Each dot represents the average success rate over 200 trials, each of which simulates 500 network measurements.

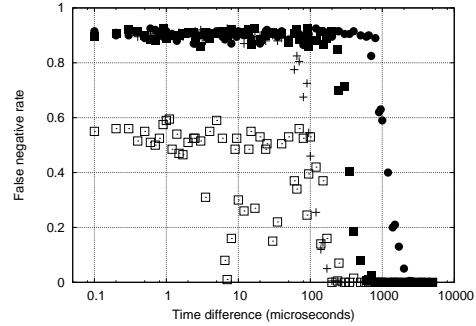


Fig. 13. Scatter plot of the time difference versus the FN rate for four dataset W_D hosts. Each dot represents the average success rate over 200 trials, each of which simulates 500 network measurements.

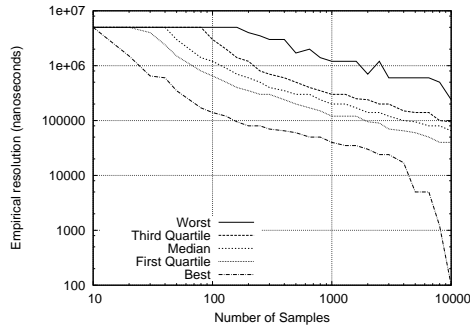


Fig. 14. Empirical resolution on dataset W_D as a function of the sample count.

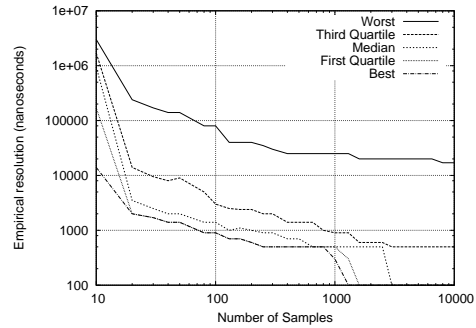


Fig. 15. Empirical resolution on the local network (Dataset L) as a function of the sample count.

and 100th percentile), and plot those with respect to the simulated measurement count in Figure 14.

As expected, increasing the number of queries radically improves the empirical resolution. Where the best dataset W_D host, with 500 measurements, could resolve $\Delta = 50\mu\text{s}$ with $< 5\%$ FP and FN rate, with 10,000 queries 31 hosts could resolve $\Delta < 35\mu\text{s}$ with a $< 5\%$ FP and FN rate. Empirical resolution improves for all other hosts as well. We can also see that the network location of an attacker can make a significant difference. The top 25% hosts can resolve a difference 3 times smaller than the bottom 25% hosts, regardless of the number of samples. This implies that an Internet attacker will do much better if many machines are available, allowing the best positioned machine to be used for attack measurements.

Figure 15 shows the same simulation performed over our 8 LAN hosts. Unfortunately, we did not sample enough discrete processing times below 100ns to directly compute the empirical resolution for these small time steps. However, we can observe a similar improvement in the empirical resolution of LAN hosts as they measure more timing samples. It is reasonable to expect that a LAN attacker making 10,000 queries might accurately

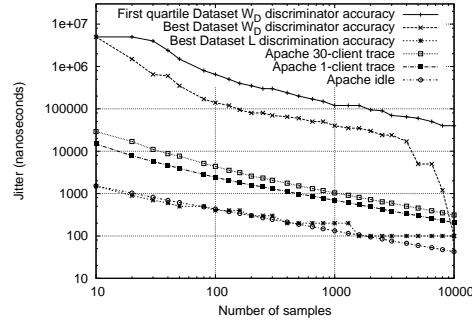


Fig. 16. Unfiltered application jitter induced by a loaded Apache web server compared to empirical resolution from our local (L) and PlanetLab (D) measurements.

resolve $\Delta < 50\text{ns}$.

6.5 Attack applications

While most of this paper has focused abstractly on the empirical resolution of an attacker to discriminate remote processing times, it is important to consider the processing times of common operations that an attacker might wish to attack.

Brumley and Boneh’s attack on SSL exploits a timing difference that is a function of the RSA private key and the message being exponentiated [Brumley and Boneh 2004]. Their attack creates two sets of inputs that *on average* have a difference. Depending on the particular input taken from those sets, the actual difference will vary. Many samples may need to be averaged together. As a result, our techniques do not directly let us infer the feasibility of their attack, but can be used to identify effective filtering strategies.

We tested their attack and found that on a 1.8GHz Athlon, it depends on timing differences of about 30,000 clock cycles ($16\mu\text{s}$). From Figure 15, we can estimate that the attack should succeed on any host within our LAN with a few hundred measurements. By plotting a comparable histogram (omitted here), we estimate that 4 hosts in dataset W_D could perform the attack with fewer than 10,000 measurements for each bit of RSA key, while others would require significantly more. Brumley and Boneh found their attack did not work over the Internet, but our results suggest that an attacker with access to enough machines, across the Internet, might well have one that works.

6.6 Application jitter

Our analysis of network jitter used a custom application on an unloaded server. Servers, like the network, can introduce latency when under load. A realistic application might introduce its own jitter into the response time because of cache misses, page faults, TLB faults, lock contention and other sources. To characterize such jitter and the ability to filter it, we analyzed the Apache web server, version 1.3.33, both under no load and under two trace-driven loads with a 1.1GB working set (larger than system RAM, guaranteeing pressure on the disk system). For simplicity, the load generation clients ran on the same computer as the HTTP server. One run was performed with sequential requests and another run was performed with 30 concurrent requests. The CPU was saturated during both runs.

Measurements were taken by a custom HTTP client, on a separate system, requesting

a small file over an unloaded LAN. We collected 100,000 measurements over the three different workloads. As an attacker would do, our client avoids connection startup jitter by first sending only part of the HTTP request, then deliberately delaying to ensure that it is received and processed by the server before it sends the remaining part of the request.

In Figure 16, we graph the unfiltered application jitter as a function of the number of request samples. We estimate the application jitter as the standard deviation of the lowest variance quantile filter, with the variance calculated by choosing 10,000 random subsets of n measurements, filtering them using a quantile filter q_i , and computing the variance of the empirical quantile \hat{q}_i across all 10k subsets.

We also plot quartiles of empirical resolution from Figures 14 and 15 from Dataset L and D trials for comparison. We can see that a realistic application like Apache introduces only a microsecond of jitter at 1000 samples even when highly loaded, and less jitter than the local network when unloaded.

There are a wide variety of different ways to architect a server. Our goal examining in Apache is to identify how effectively the jitter could be filtered from one sample real world application. Of course, different application structures, on different operating systems, and on different hardware may yield more or less jitter.

7. RELATED WORK

Side channel attacks. The timing attacks discussed in this paper are an example of side channel attacks, where a system leaks information due to its physical implementation. An early example of such an attack is the password authentication weakness discovered in the Tenex operating system [Lampson 1983]. Kocher was the first to observe that side-channels attacks could generally be applied against common cryptographic algorithms. His analysis using device response times [Kocher 1996] and power consumption [Kocher et al. 1999] to derive cryptographic secrets were the basis for much of the recent work in this field. In the recent Advanced Encryption Standard (AES) competition, the ciphers were examined for the potential of side-channel attacks [Daemen and Rijmen 1999].

Side channels exist where a computer may leak its internal state through RF emissions [Kuhn and Anderson 1998]. Strikingly, methods as effortless as watching reflections of displays, at a distance, may allow an observer to see the screen [Kuhn 2002; Backes et al. 2008]. Side channels have been used to detect passwords over SSH, through the use of keystroke timing [Song et al. 2001]. Aciçmez and Çetin Kaya Koç [2006] presents an attack on AES that exploits a side channel when an attacker can see the stream of individual cache hits and misses.

Timing attacks have been applied to cryptosystems. Kelsey et al. [2000] conjectured that cache miss behavior may also be used as a side channel. Page [2002] later presented an analysis of caching behavior and described an attack against DES. Bernstein [2005] showed that AES's CPU cache-miss behavior leaks key bits. A more powerful attack on AES exploiting shared cache state has been done by Osvik et al. [2006]. Silverman and Whyte [2007] summarize a timing attack on the NTRU cryptosystem that exploits a timing difference on the number of SHA-1 computations. Schindler [2000] described an attack against the chinese remainder theorem in RSA and Schindler [2002] modeled and optimized attacks against RSA. The Brumley and Boneh [2004] attack on RSA was improved by Aciçmez et al. [2005].

More recently, work has been done showing a side channel from a shared branch predic-

tor or trace cache that leaks the execution path taken by a program and permitting attacks on cryptosystems [Aciçmez et al. 007c]. Aciçmez et al. [007a] show a local attack on RSA that exploit branch mispredictions delays in order to determine the secret key. Aciçmez et al. [007b] present an interprocess timing attack over the loopback interface on AES that exploits CPU cache timing differences using about 10^{11} encryptions and were able to distinguish the correct AES key among 2^{12} alternatives.

Timing attacks have been applied to the Internet. Kohno et al. [2005] showed that it is possible to fingerprint a host on the Internet by using TCP or ICMP timestamps to measure differences between machine clock skews as tiny as $1\mu\text{s}$ per second. Felten and Schneider [2000] shows how servers can fingerprint anonymous web clients through detecting the timing difference between a client cache hits and misses. And, Bortz et al. [2007] showed that timing difference may leak secrets such as the existence of an account or shopping cart size in web applications.

Network measurement. There have been many attempts to characterize the end-to-end behavior of the Internet. The most comprehensive work is Paxson's measurements and analysis of Internet end-to-end dynamics [Paxson 1997a; 1997b]. Paxson characterized such issues as routing pathologies (e.g., routing loops), outages, flutter, and the stability and lifetime of routing information. He also examined Internet packet dynamics, including the effects of packet loss, corruption, and reordering on TCP.

The earliest studies of network latency and jitter focused on these attributes because of their effect on important parameters in the TCP protocol. An accurate value for the round trip time is needed to estimate the correct values for TCP retransmit timers and window size [Claffy et al. 1993]. If jitter causes the round trip time to be incorrectly measured, the TCP protocol may incorrectly initialize its timers.

Internet path delay time was characterized as a shifted gamma distribution in Mukherjee [1994]. His measurements used standard ICMP echo requests, and achieved millisecond precision. Many other researchers have performed end-to-end assessments of Internet packet behavior [Acharya and Saltz 1996; Bolot 1993; Sanghi et al. 1993; Wang et al. 2003]. Barford and Crovella [2001] characterized causes of delay and jitter in a web server scenario. Casner et al. [2001] measured Internet jitter to $20\mu\text{s}$ resolution on a wide area backbone network to study the feasibility of using backbone IP networks as virtual circuits.

Generally, these studies were concerned with millisecond-scale events, and did not consider the notion of an attacker willing to make thousands or even millions of repeated queries in order to gain increased timing accuracy of a remote machine's processing time.

Clock synchronization. Clock synchronization and remote timing attacks must both handle Internet jitter and delay. Unlike clock synchronization algorithms, an attacker only needs to worry about the stability of one clock — their own, and over a timescale of minutes. The attacker can also afford to collect many measurements.

The Network Time Protocol [Mills 1991; 1992] is designed to synchronize the system clocks of computers over the network. NTP must, by necessity, measure and compensate for network latency and jitter, but its goals are to achieve millisecond, not sub-microsecond accuracy.

Many types of network measurement depend more upon low clock skew variation across the measurements hosts than offsets from real time [Paxson 1998]. Protocols other than NTP have been designed to minimize clock skew [Pasztor and Veitch 2002; Veitch et al.

2004].

8. CONCLUSION

This paper studied the scope and potential of performing a remote timing attack, both on a LAN and across the Internet. Any such attack will require making multiple timing measurements of an event on a remote server and filtering those measurements to eliminate noise that might be induced by the network or by the end hosts. We have shown that, even though the Internet induces significant timing jitter, we can reliably distinguish remote timing differences as low as $20\mu\text{s}$. A LAN environment has lower timing jitter, allowing us to reliably distinguish remote timing differences as small as 100ns (possibly even smaller). These precise timing differences can be distinguished with only hundreds or possibly thousands of measurements.

Good filtering of those measurements is fundamental to mounting a successful attack. Contrary to conventional wisdom, using either the median response time or the minimum response time observed as a filter significantly under-performs filters that sort the data and look at values early in the range (e.g., 1% into the sorted list). Based on filters that use these low percentiles, we can construct a “box test” that can reliably distinguish small timing differences, when they are present, with low false positive and low false negative rates.

We also observed, generally, that the round trip time or network hop count did not significantly contribute to the network jitter, and thus network distance may not confer immunity to remote timing attacks. We found that the choice CPU or networking card may introduce more jitter than a local area network. Prospective attackers can work around this by benchmarking their measurement machines, in advance.

If an attacker can accurately perform timing measurements, then a number of cryptographic or algorithmic weaknesses in a server might leak critical information to the attacker. As a consequence, we recommend that the algorithms used inside web and other Internet servers that process important secrets be carefully audited and, where necessary, be modified to limit observable differences in execution times to at most a few microseconds.

REFERENCES

- ACHARYA, A. AND SALTZ, J. 1996. A study of Internet round-trip delay. Tech. Rep. CS-TR-3736, Department of Computer Science, University of Maryland. Dec.
- ACIİÇMEZ, O. AND ÇETIN KAYA KOÇ. 2006. Trace-driven cache attacks on aes (short paper). In *Information and Communications Security, 8th International Conference*. Raleigh, NC, 112–121.
- ACIİÇMEZ, O., ÇETIN KAYA KOÇ, AND SEIFERT, J.-P. 2007(a). Predicting secret keys via branch prediction. In *The Cryptographers' Track at the RSA Conf. (CT-RSA)*. San Francisco, 225–242.
- ACIİÇMEZ, O., SCHINDLER, W., AND ÇETIN K. KOÇ. 2005. Improving brumley and boneh timing attack on unprotected ssl implementations. In *Proceedings of the 12th ACM conference on Computer and communications security*. Alexandria, VA, USA, 139–146.
- ACIİÇMEZ, O., SCHINDLER, W., AND ÇETIN KAYA KOÇ. 2007(b). Cache based remote timing attack on the AES. In *The Cryptographers' Track at the RSA Conf. (CT-RSA)*. San Francisco, 271–286.
- ACIİÇMEZ, O., SEIFERT, J.-P., AND ÇETIN KAYA KOÇ. 2007(c). Micro-architectural cryptanalysis. *IEEE Security and Privacy* 5, 4, 62–64.
- BACKES, M., DÜRUMUTH, M., AND UNRUH, D. 2008. Compromising reflections or how to read lcd monitors around the corner. In *IEEE Symposium on Security and Privacy*. Oakland, CA.
- BARFORD, P. AND CROVELLA, M. 2001. Critical path analysis of TCP transactions. *IEEE/ACM Trans. Netw.* 9, 3, 238–248.

- BERNSTEIN, D. J. 2005. Cache-timing attacks on AES. <http://cr.yp.to/papers.html#cachetiming>.
- BOLOT, J.-C. 1993. End-to-end packet delay and loss behavior in the Internet. In *ACM SIGCOMM: Conf. on Applications, Technologies, Architectures, and Protocols for Computer Communication*. San Francisco, CA, 289–298.
- BORTZ, A., BONEH, D., AND NANDY, P. 2007. Exposing private information by timing web applications. In *16th Int. World Wide Web Conf.* Banff, Alberta, Canada.
- BRUMLEY, D. AND BONEH, D. 2004. Remote timing attacks are practical. In *Proc. of the 12th USENIX Security Symposium*. Washington, DC.
- CANVEL, B., HILTGEN, A., VUAGNOUX, M., AND VAUDENAY, S. 2003. Password interception in a TLS/SSL channel. In *Advances in Cryptology — Crypto Proc.* LNCS, vol. 2729. Springer-Verlag, 583–599.
- CARTER, J. L. AND WEGMAN, M. N. 1979. Universal classes of hash functions. *J. Comput. Syst. Sci.* 18, 2 (Apr.), 143–154.
- CASNER, S., ALAETTINOGLU, C., AND KUAN, C.-C. 2001. A fine-grained view of high-performance networking. NANOG22 meeting. <http://www.nanog.org/mtg-0105/casner.html>.
- CHUN, B., CULLER, D., ROSCOE, T., BAVIER, A., PETERSON, L., WAWRZONIAK, M., AND BOWMAN, M. 2003. PlanetLab: An overlay testbed for broad-coverage services. *ACM Computer Communications Review* 33, 3 (July).
- CLAFFY, K. C., POLYZOS, G. C., AND BRAUN, H.-W. 1993. Measurement considerations for assessing unidirectional latencies. *Internetworking: Research and Experience* 4, 3 (Sept.), 121–132.
- CROSBY, S. AND WALLACH, D. S. 2003. Denial of service via algorithmic complexity attacks. In *Proc. of the 12th USENIX Security Symposium*.
- DAEMEN, J. AND RIJMEN, V. 1999. Resistance against implementation attacks: A comparative study of the AES proposals. In *2nd AES Candidate Conf.* Rome, Italy.
- FELTEN, E. W. AND SCHNEIDER, M. A. 2000. Timing attacks on web privacy. In *Proc. of 7th ACM Conf. on Computer and Communications Security*. Athens, Greece.
- JENKINS, R. J. 1995. Hash functions for hash table lookup. <http://burtleburtle.net/bob/hash/evahash.html>.
- KELSEY, J., SCHNEIER, B., WAGNER, D., AND HALL, C. 2000. Side channel cryptanalysis of product ciphers. *Journal of Computer Security* 8, 2-3, 141–158.
- KOCHER, P. 1996. Timing attacks on implementations of Diffie-Hellman, RSA, DSS, and other systems. In *Advances in Cryptology — Crypto Proc.*, N. Koblitz, Ed. LNCS, vol. 1109. Springer-Verlag, Santa Barbara, California.
- KOCHER, P., JAFFE, J., AND JUN, B. 1999. Differential power analysis. In *Advances in Cryptology — Crypto Proc.*, M. Wiener, Ed. LNCS, vol. 1666. Springer-Verlag, Santa Barbara, CA.
- KOHNO, T., BROIDO, A., AND CLAFFY, K. 2005. Remote physical device fingerprinting. In *Proc. of the IEEE Symposium on Security and Privacy*. Oakland, CA.
- KUHN, M. G. 2002. Optical time-domain eavesdropping risks of CRT displays. In *Proc. of the IEEE Symp. on Security and Privacy*. Oakland, CA.
- KUHN, M. G. AND ANDERSON, R. J. 1998. Soft Tempest: Hidden data transmission using electromagnetic emanations. In *2nd Workshop on Information Hiding*. Portland, OR.
- LAMPSON, B. W. 1983. Hints for computer system design. *ACM Operating Systems Review* 15, 5 (Oct.), 33–48. Reprinted in *IEEE Software* 1(1):11-28 (Jan. 1984).
- MILLS, D. L. 1991. Internet time synchronization: The Network Time Protocol. *IEEE Transactions on Communications* 39, 10 (Oct.), 1482–1493.
- MILLS, D. L. 1992. Network time protocol (version 3). Tech. Rep. RFC-1305, Internet Engineering Task Force. Mar. <ftp://ftp.rfc-editor.org/in-notes/rfc1305.txt>.
- MOORE, D. AND MCCABE, G. 2004. *Introduction to the Practice of Statistics*, 5th ed. Freeman, New York.
- MUKHERJEE, A. 1994. On the dynamics and significance of low frequency components of Internet load. *Internetworking: Research and Experience* 5, 4, 163–205.
- OSVIK, D. A., SHAMIR, A., AND TROMER, E. 2006. Cache attacks and countermeasures: The case of AES. In *The Cryptographers' Track at the RSA Conf. (CT-RSA)*. LNCS, vol. 3860. Springer, San Jose, CA, 1–20.

- PAGE, D. 2002. Theoretical use of cache memory as a cryptanalytic side-channel. Tech. Rep. CSTR-02-003, Department of Computer Science, University of Bristol. June.
- PASZTOR, A. AND VEITCH, D. 2002. PC based precision timing without GPS. In *Proc. of the ACM SIGMETRICS Int. Conf. on Measurement and Modeling of Computer Systems*. Marina Del Rey, CA, 1–10.
- PAXSON, V. 1997a. End-to-end Internet packet dynamics. In *ACM SIGCOMM: Applications, Technologies, Architectures and Protocols for Computer Communication*. Cannes, France.
- PAXSON, V. 1997b. Measurements and analysis of end-to-end Internet dynamics. Ph.D. thesis, University of CA at Berkeley, Berkeley, CA.
- PAXSON, V. 1998. On calibrating measurements of packet transit times. In *SIGMETRICS/PERFORMANCE: Joint Int. Conf. on Measurement and Modeling of Computer Systems*. Madison, WI, 11–21.
- SANGHI, D., AGRAWALA, A. K., GUDMUNDSSON, O., AND JAIN, B. N. 1993. Experimental assessment of end-to-end behavior on Internet. In *Twelfth Annual Joint Conf. of the IEEE Computer and Communications Societies*. San Francisco, CA, 867–874.
- SCHINDLER, W. 2000. A timing attack against RSA with the chinese remainder theorem. In *Cryptographic Hardware and Embedded Systems*. Worcester, MA, 109–124.
- SCHINDLER, W. 2002. Optimized timing attacks against public key cryptosystems. *Statistics and Decisions* 20, 191–210.
- SILVERMAN, J. H. AND WHYTE, W. 2007. Timing attacks on NTRUEncrypt based on variation in number of hash calls. In *The Cryptographers' Track at the RSA Conf. (CT-RSA)*. San Francisco.
- SONG, D. X., WAGNER, D., AND TIAN, X. 2001. Timing analysis of keystrokes and timing attacks on SSH. In *Proc. of the Tenth USENIX Security Symposium*. Washington, DC.
- VEITCH, D., BABU, S., AND PASZTOR, A. 2004. Robust synchronization of software clocks across the Internet. In *Proc. of the 4th ACM SIGCOMM Conf. on Internet Measurement*. Taormina, Sicily, Italy, 219–232.
- WANG, Z., ZEITOUN, A., AND JAMIN, S. 2003. Challenges and lessons learned in measuring path RTT for proximity-based applications. In *Passive and Active Measurement Workshop*. San Diego, CA.

Received May 2007; revised May 2008; accepted Month Year