

RESEARCH ARTICLE

Nonintrusive tracing in the Internet

Alina Olteanu¹, Yang Xiao^{2*}, Jing Liu², Thomas M. Chen³ and C. L. Philip Chen⁴¹ Spiru Haret University, Bucharest, Romania² Department of Computer Science, The University of Alabama, Tuscaloosa, AL 35487, U.S.A.³ School of Engineering, Swansea University, Swansea, Wales SA2 8PP, U.K.⁴ Faculty of Science of Technology, University of Macau, Macau, China

SUMMARY

Intruders that log in through a series of machines when conducting an attack are hard to trace because of the complex architecture of the Internet. The thumbprinting method provides an efficient way of tracing such intruders by determining whether two connections are part of the same connection chain. Because many connections are transient and therefore short in length, choosing the best time interval to thumbprint over can be an issue. In this paper, we provide a way to shorten the time interval used for thumbprinting. We then study some special properties of the thumbprinting function. We also study another mechanism for tracing intruders in the Internet based on a timestamping approach, which passively monitors flows between source and destination pairs. Given a potentially suspicious source, we identify its true destination. We compute the error probability of our algorithm and show that its value decreases exponentially as the observation time increases. Our simulation results show that our approach performs well. Copyright © 2012 John Wiley & Sons, Ltd.

KEYWORDS

nonintrusive; tracing; Internet

*Correspondence

Yang Xiao, Department of Computer Science, The University of Alabama, Tuscaloosa, AL 35487, U.S.A.

E-mail: yangxiao@ieee.org

1. INTRODUCTION

Constant change is perhaps a guiding principle in the Internet. Recent advances in technology have led to significant growth of the Internet by factors of 10^3 and 10^6 in backbone speed and the number of hosts, respectively [1]. Since the public expansion of the Internet in 1990, many new challenges have surfaced; among them, operations between untrusted end points, more demanding applications, and less sophisticated users have caused severe stress on the Internet requirements. Furthermore, the number of attacks on networked computer systems has been growing exponentially from year to year. When considering the task of tracing intruders in the Internet, three main challenges must be taken into account. First, attackers hide their origin by making use of the Internet's architecture. By using different hosts that belong to different countries and administrative domains to route malicious acts, intruders' actions become extremely difficult to trace. Second, the data collected from an Internet trace is usually incomplete or has missing values. For example, different domains of Internet service providers may not share data because of access issues, such as when they are in different

countries. Finally, routes change frequently, packets are lost, and the network latency and time to convergence can be significantly increased because of routing instability.

To deal with such problems, two tracing mechanisms exist [2]: (i) methods of keeping track of all individuals and accounting for all activities and (ii) reactive tracing, in which no global accounting is attempted until a problem arises. In our view, the first mechanism is most related to network accountability, and the second mechanism is most related to network forensics. Network-based tracing and host-based tracing are two common approaches for reactive tracing in connection chains. Host-based tracing involves one tracing system per network host [2], and a chain of connection hosts can be known via (i) host communications [3] or (ii) reversing the attack chain by breaking the hosts in reverse order [4]. Host-based tracing schemes suffer when an extended connection crosses a host not running the system [2]. Network-based tracing has the advantage of not relying on hosts that can be untrustworthy or requiring host participation. Instead, network-based tracing uses the invariance of connections at higher protocol layers (such as the transport layer) to establish whether two connections are part of the same connection chain.

We study two approaches for network-based tracing. The first approach is based on the idea of a connection thumbprint. A thumbprint is very similar to a checksum or the computed summary of the content of a connection [2,5]. The thumbprint summarizes the content of a connection by using only a small amount of data while simultaneously preserving the uniqueness of the connection. Therefore, thumbprints of related connections are similar and can lead to the construction of a connection chain. However, the thumbprinting method is dependent on the duration of the connection, so the time used for thumbprinting is very important especially in the case of transient flows. In this paper, we provide a way to shorten the time interval required for thumbprinting. We study the tradeoff introduced by using a smaller time interval and its effect on basic thumbprint properties such as sensitivity and robustness.

Because thumbprinting relies on the content of the packages and therefore can be counterattacked by encryption, we study a second mechanism for tracing Internet intruders that involves passively monitoring flows between source and destination pairs in the Internet, a process similar to [6]. This approach is based on monitoring the transmission activities of nodes and does not interfere with network operations. A one-hop communication graph can be constructed by matching transmission timestamps and acknowledgements. On the basis of this graph, the nodes that are part of a connection chain are those which communicate at a sufficiently high rate. Our approach is different from the approach in [6] in that we make fewer statistical assumptions and we use clearer, simpler derivations. In addition, because we are working with the Internet, we face other challenges, such as the availability of only partial information and the route instability mentioned earlier. To account for this, we assume that the information on every link has a certain probability of availability.

Our motivations in this paper can be derived from the previous discussions as follows:

- (1) To find a way to improve current thumbprinting method for transient traffic flows by shortening the time interval required for thumbprinting.
- (2) To study the tradeoff introduced by using a smaller time interval and its effect on basic thumbprint properties such as sensitivity and robustness.

According to these motivations, relevant contributions are made as follows:

- (1) To find out the equation for sampling traffic in any time interval and the properties that can be used for identifying related flows so that we can shorten the time interval required for thumbprinting.
- (2) To obtain the error probability of our algorithm and show that its value decreases exponentially as the observation time increases.

- (3) To consider network background noise and redesign a passive tracing scheme accordingly.
- (4) To provide numerical results and simulation results.

The remainder of the paper is organized as follows. Section 2 presents significant work related to our problem, Section 3 presents our results concerning minimizing the thumbprinting interval, Section 4 contains a detailed introduction to the tracing algorithm followed by an evaluation of its consistency in terms of error probability, Section 5 presents some graphical interpretations of the derivations, Section 6 provides simulations, and Section 7 is the paper's conclusion.

2. RELATED WORK

Paper [2] presents an Internet Protocol traceback method based on the thumbprint idea. Often, attacks are conducted by logging through a chain of host machines. In this way, intruders make use of the Internet architecture to hide their true origin. The thumbprint technology, which is the summary of a connection, is used to compare connections and therefore trace the source of such an attack. A thumbprint is computed for each time interval, typically a minute, of each connection [2]. Thumbprinting has the advantage of preserving the characteristics of a connection while being cheap to compute and requiring little storage. Similar methods such as checksum and message digest have two main disadvantages. Firstly, an error in the content leads to a different value being computed, and thus these methods are not robust. Secondly, they are not additive (i.e., two successive values cannot be combined to provide a new value for a longer interval). Other compression schemes require far more storage space than thumbprinting.

As the thumbprinting method is dependent on the duration of the connection, it is of importance to shorten the time interval used for thumbprinting, especially in the case of transient flows. This is the purpose of our derivation in Section 3.

However, fingerprinting techniques are vulnerable and may not be of interest when the attacker uses encryption at every hop in the chain. To deal with such problems, techniques based on packet timing information are more suitable.

Paper [6] presents a mechanism for tracing intruders in anonymous mobile ad hoc networks (MANETs) by passively monitoring flows between source and destination pairs. The algorithm is concerned with tracing the destination of a certain source that is considered to be suspicious according to a higher layer intrusion detection protocol. Because anonymity of nodes in the MANET is assumed and the monitoring is performed in a nonintrusive way, this research only considers the timestamps of packets on adjacent links and uses these to establish causal relationships between packets. A graph is used to model the communications among nodes. On the basis of traffic analysis and information carried in the packets, the graph

is partitioned into two parts, one containing the set of possible destinations and one with which the source does not communicate. The following assumptions are made about the distribution of traffic: the transmission activities on each link are assumed to follow a Poisson process, the rate of the flow to be traced is sufficiently high, and the duration of the observation is long enough. Therefore, the method will not work for low-rate or transient flows.

Much like [6], we focus on tracing flows between pairs (source and destination), but we do so with the Internet. Unlike [6], we assume differences in data due to propagation delays, clocks being skewed during observation at multiple nodes at the same time, and the data that is either missing from observations or observed erroneously. We account for this by setting a probability at which information is available on every link. We also make fewer assumptions about the rate of the Poisson process, and the calculations are significantly simplified. Some other related work includes [7–41].

3. THUMBPRINTING

The thumbprinting approach is based on the fact that, at higher protocol layers, the content of a connection is invariant at all points of the chain. The thumbprint summarizes the content of the connection by using only a small amount of data while preserving the uniqueness of the connection. We first define the terminology and then present a method for shortening the time interval required for thumbprinting and some interesting properties of the thumbprint function.

3.1. Original algorithm

A thumbprint is a function of a connection that preserves the unique characteristics of the connection. We use the notations and definition of a thumbprint from [2]. Consider a sequence of transmitted characters a_1, a_2, \dots, a_L that needs to be thumbprinted. Consider also the function $\alpha: A \rightarrow R^K$, which takes a character and returns a vector of K real numbers. Let $\alpha_k(\cdot)$ denote the k th component of the function. Here, K is a short fixed number representing the number of thumbprint components. If we consider L as the period of a connection, we can associate frequencies $\bar{f} = (f_1, f_2, \dots, f_L)$ to our character sequence. Then, the thumbprint is defined as [2]:

$$T_k = \sum_{i=1}^L \alpha_k(a_i) f_i, \quad k \in \{1, 2, \dots, K\}. \quad (1)$$

The thumbprint is thus a linear combination of the frequencies of characters and their corresponding weights. T_k represents the k th component of the K -dimensional thumbprint vector.

Each thumbprint component, $T_k(C, t)$, is a function of a specific connection C and the time interval t , in which the

thumbprints have been computed. From [2], we also know that the comparison of two sets of thumbprints, $T_k(C, t)$ and $T_k(C', t)$, $k \in \{1, 2, \dots, K\}$, for two different connections C and C' , is given by the following formula:

$$\delta_t(C, C') = \log \left(\prod_{k=1}^K |T_k(C', t) - T_k(C, t)| \right) \quad (2)$$

δ_t represents the logarithm of the product of component differences for two thumbprints in a specific time interval t . A large absolute value of δ_t implies that the two connections are related, whereas a small absolute value suggests that they are independent.

3.2. Minimizing the thumbprinting interval

In this section, on the basis of the results in Section 3.1 [2], we try to find the best time T , over which to thumbprint. The length of time divisions for the experiments in [2] is 1 time unit. However, many connections in the Internet are transient and therefore short in length; in these cases, a shortened thumbprinting interval is needed. Hence, we need to find a new equation, rather than (1), to better describe the characteristic of each link or session in terms of more precise time duration.

We start with the thumbprint function given by (1), and we then account for the time in which the thumbprint has been computed. This gives us the following function:

$$\sum_{a,t} \alpha_k(a, t) f(a, t) \quad (3)$$

We will use the function given by (3) earlier as our thumbprinting function.

We now show the results of changes to variables $a = Ly$ and $t = T\tau$. This transforms our time interval from $[0, T]$ to $[0, 1]$ and the samples from every time unit to every $1/T$ time units. Thus, the thumbprinting function in (3) can be approximated by the following integral:

$$\iint_{[0,L] \times [0,T]} f(a, t) \alpha_k(a, t) da dt. \quad (4)$$

We have

$$\begin{aligned} \iint_{[0,L] \times [0,T]} f(a, t) \alpha_k(a, t) da dt &\approx LT \sum_{a,t} \alpha_k(a, t) f(a, t) / (LT) \\ &\stackrel{a=Ly}{t=T\tau} \\ &= LT \sum_{y,\tau} \alpha_k(Ly, T\tau) f(Ly, T\tau) / (LT) \\ &= LT \iint_{[0,1] \times [0,1]} f(Ly, T\tau) \alpha_k(Ly, T\tau) dy d\tau \end{aligned}$$

The time has become interval $[0, 1]$.

Because thumbprinting over a one unit time interval may not provide the best results, we further provide a way to transform interval $[0, 1]$ into an arbitrary time

interval $[p, q]$. In the following calculus, we use $y = (u - m)/(n - m)$ and $\tau = (v - p)/(q - p)$. We have

$$\begin{aligned}
 & LT \iint_{[0,1] \times [0,1]} f(Ly, T\tau) \alpha_k(Ly, T\tau) dy d\tau \\
 &= LT \iint_{[m,n] \times [p,q]} f(L(u - m)/(n - m), T(v - p)/(q - p)) \\
 &\quad \alpha_k(L(u - m)/(n - m), T(v - p)/(q - p)) \left| \frac{\partial(y, \tau)}{\partial(u, v)} \right| dudv \\
 &= LT \iint_{[m,n] \times [p,q]} \frac{f(L(u - m)/(n - m), T(v - p)/(q - p))}{(n - m)(q - p)} \\
 &\quad \alpha_k(L(u - m)/(n - m), T(v - p)/(q - p)) dudv.
 \end{aligned} \tag{5}$$

In Equation (5), the time has become interval $[p, q]$. This interval can be suitably chosen to accommodate connections with low/high data rates and transient connections to give the best performance possible under different scenarios.

3.3. Properties of the thumbprinting function

In the following, we describe a way to easily compute the value of the integral given by (4) by using the mean value theorem [42]. We further study some special properties of the thumbprinting function.

Consider the right-hand side of Equation (5).

Let us choose u_0 and v_0 to be the points toward which we concentrate intervals $[m, n]$ and $[p, q]$, respectively. By making $m, n \rightarrow u_0$ and $p, q \rightarrow v_0$ (i.e., $n - m \rightarrow 0$ and $q - p \rightarrow 0$), we have

$$\begin{aligned}
 & \iint_{[0,L] \times [0,T]} f(a, t) \alpha_k(a, t) da dt = LT \iint_{[0,1] \times [0,1]} f(Ly, T\tau) \alpha_k(Ly, T\tau) dy d\tau \\
 &= \frac{LT \iint_{[m,n] \times [p,q]} f\left(L \frac{u - m}{n - m}, T \frac{v - p}{q - p}\right) \alpha_k\left(L \frac{u - m}{n - m}, T \frac{v - p}{q - p}\right) dudv}{(n - m)(q - p)} \\
 &= LT f\left(L \frac{u_0 - m}{n - m}, T \frac{v_0 - p}{q - p}\right) \alpha_k\left(L \frac{u_0 - m}{n - m}, T \frac{v_0 - p}{q - p}\right) \\
 &= LT h\left(L \frac{u_0 - m}{n - m}, T \frac{v_0 - p}{q - p}\right).
 \end{aligned} \tag{6}$$

We have considered $h = f\alpha_k$, where h is a continuous function. Let $I = LT \iint_{[0,1] \times [0,1]} f(Ly, T\tau) \alpha_k(Ly, T\tau) dy d\tau$. If h is a C^n function, where $n \geq 1$, then the set of points $\gamma = \{(u_0, v_0); I = LT(f\alpha_k)(u_0, v_0)\}$ represents a C^n curve that, according to the Implicit Function theorem [43], ensures the existence of a function φ such that $v_0 = \varphi(u_0)$, $u_0 \in J$, where J is an interval included in $[m, n]$. This way, we can express one of the variables u_0 and v_0 as a function of the other.

The last equality in (6) is true for any $(u_0, v_0) \in [m, n] \times [p, q]$ on the curve γ .

We have thus shown that the value of the integral in (4) can be found by computing the value of the function $f\alpha_k$ at any point belonging to curve γ .

Furthermore, if we choose

$$\begin{aligned}
 n &= m + \varepsilon, \quad u_0 = m + \varepsilon t, \quad 0 < t < 1, \quad q = p + \delta \\
 \varphi(u_0) &= p + \lambda(t)\delta
 \end{aligned}$$

we have

$$\begin{aligned}
 I &= LT(f\alpha_k)\left(L \frac{\varepsilon t}{\varepsilon}, T \frac{\lambda(t)\delta}{\delta}\right) \\
 &= LT f(Lt, T\lambda(t)) \alpha_k(Lt, T\lambda(t)), \quad t \in J
 \end{aligned}$$

We have shown that on the curve $\gamma = \{(u_0, \varphi(u_0)), u_0 \in J\}$, we have $\alpha_k(u_0, \varphi(u_0)) = \frac{I}{LT f(u_0, \varphi(u_0))} = \text{const}$. To determine the constant, we simply need to compute I . Therefore, α_k and f are inversely proportional on curve γ .

We now introduce an abstraction of the derivations used earlier by using a mean value operator. The use of this operator generalizes the procedure of associating the function h with its mean value. Such an operator can be defined as follows:

$$T(h)(0) = h(0) = \lim_{x \downarrow 0} \frac{\int_0^x h(t) dt}{x} = \lim_{x \downarrow 0} T(h)(x), \quad h \in C([0, b]) \text{ for one variable and}$$

$$\begin{aligned}
 T(h)(0, y) &= \lim_{x \downarrow 0} \frac{\iint_{[0,x] \times [0,y]} h(t_1, t_2) dt_1 dt_2}{xy} \\
 &= \frac{\int_0^y h(0, t_2) dt_2}{y}, \quad y \neq 0
 \end{aligned}$$

$$T(h)(x, 0) = \frac{\int_0^x h(t_1, 0) dt_1}{x}$$

$T(h)(0, 0) = h(0, 0)$ for two variables.

Operator T associates function h with another function that is equal at every point x to the mean value of function h on interval $[0, x]$. In Section 5, we will see that the computation has been simplified by using this mean value approach.

4. TRACING THE TRUE DESTINATION OF A SOURCE

Because many attackers use encryption at every hop in the chain, the same packet will have different content on every link, and fingerprinting techniques are not appropriate in this case. Therefore, we present a mechanism from paper [6] for tracing intruders in anonymous MANETs that are based on using transmission timestamps to passively monitor the flows between source and destination pairs. Such an

approach can be implemented by using sensors equipped with energy detectors to measure transmission timestamps and then distributing them over the field of interest. All measurements are fused and processed by a centralized monitor at a fusion center.

We first introduce the algorithm and its main features and then analyze the algorithm in terms of the error probability of detecting the true destination of a source. We show that the error probability decays exponentially with the observation time and also prove a similar result when the number of observed graphs goes to infinity.

4.1. Original algorithm

Our derivation is based on a variation of the tracing algorithm from [6]. In [6], the approach for tracing the destination of a source is twofold. First, on the basis of the transmission activities on adjacent links, it can be determined whether these links are part of the same flow. Then, a set of possible destinations is computed. Second, the intersection method is used: using the changes in topology, some nodes are eliminated, leading to a smaller set of possible destinations.

It is assumed that the transmission activities on each link follow a Poisson process S , and its realization is denoted by s . Therefore, a realization of the transmission timestamps of data packets on link 1 will be $s_1 = (s_1(1), s_1(2), s_1(3), \dots, K)$. Here, uppercase letters denote random variables, lowercase letters realizations, boldface letters vectors, and plain letters scalars.

In the following, we briefly introduce the idea behind the traffic analysis method and the intersection of different topologies method from [6].

For traffic analysis, consider two realizations (s_1, s_2) of the transmission activities (timestamps) on two adjacent links. Let m and n be the indices of the current timestamps on links 1 and 2, respectively ($0 \leq m \leq \delta_1, 0 \leq n \leq \delta_2$). These sets of timestamps are matched against each other sequentially by assessing the difference $s_2(n) - s_1(m)$. If $s_2(n) - s_1(m) \leq \Delta$, where Δ is a predefined maximum delay, and this difference is nonnegative, then the two timestamps match. The matching timestamps from the two connections form a pair of sequences (f_1, f_2) . Given the number of matching timestamps on a link, the empirical rate on that link can be now estimated. Specifically, if f_i contains $|F_i|$ timestamps over time T , then the empirical rate is $|F_i|/T$.

Let $\tau \geq 0$ be a given rate. Given a graph sourced at j , by repeatedly applying the matching timestamp algorithm for pairs of adjacent links and selecting only the flows that support rates of at least τ , a subgraph of the initial graph can be obtained. This new graph will contain only the nodes that node j can talk to at a rate higher or equal to τ . The method, called Trace Destination (TD), is applied for the graph sourced at O , and a set of possible destinations is obtained. However, node O can only be a relay node for some of these flows. To trace down the flows that are going through O , but did not originate at O , TD is applied for all immediate predecessors of O and the corresponding graphs sourced at these nodes. The set of

obtained destinations is then subtracted from the previously obtained set of destinations of O .

To speed up the algorithm's convergence, a changing topology feature is used. Basically, the algorithm TD is applied for every communication graph that is observed. For each observed graph, a set of destinations is obtained. The final set of possible destinations is the intersection over all sets from all observed graphs. Out of this final set, the node that appears most frequently in all topologies is selected. If there are multiple possible destinations with the same number of appearances across all different topologies, then one node is selected at random.

4.2. Analysis

Our tracing algorithm is a variation of the one from Section 4.1, in which we account for the lack of information in the observations. We model the lack of information by considering that the information on each link is available with a certain probability. We denote this probability by $p_k, k \in \{1, 2, \dots, K\}$, where K is the number of links in the considered path. In this modified algorithm, only flows that support rates of at least $r_k/p_k \geq \tau$ are selected from the initial graph.

It is now useful to study how well our variation of the algorithm from Section 4.1 converges. We thus study the probability P_e with which the algorithm finds the destination erroneously (called error probability) and analyze its asymptotic behavior over time T . We would like to mention that even though we assume Poisson distribution of packets on a link for our specific derivation, the algorithm and analysis here can be applied to different types of traffic.

Our error probability derivation is also inspired from [6]. Assume that the transmission activities on each link follow a Poisson process of rate $R < 1$. It is well known that the Poisson distribution of the number of events in a time interval $(t, t + t')$ is given by the following relation:

$$P[(m_{t+t'} - m_t) = n] = \frac{e^{-Rt'} (Rt')^n}{n!}, n = 0, 1, \dots$$

where $m_{t+t'} - m_t$ represents the number of events in time interval $(t, t + t')$ (see for example [44]).

Let O be the source whose destination we are trying to find, and let θ be the true destination of O . We denote the destination rendered by our algorithm with $\hat{\theta}$. We then derive the probability that $\hat{\theta}$ is not the true destination of O . Let n_k be the number of timestamps taken in time interval (m_k, m_{k+1}) , corresponding to link $(k, k + 1)$. We denote the total number of timestamps taken for one flow by N_1 .

By definition,

$$\begin{aligned} P_e &:= P_e(TD) = P(\theta \neq \hat{\theta}) \\ &= P(\theta \text{ is not correctly detected by algorithm TD}) \end{aligned}$$

According to [6], Theorem 4.4, the fact that the algorithm does not find the true destination is due to one of three possibilities:

- (1) θ is not detected because the empirical rate along the path from O to θ is less than τ ; we will denote the probability of this event by $P(A)$.
- (2) θ is mistaken for a relay node because there is some node j , which is a successor of θ , for which the empirical rate from O to j is greater than τ . We denote the probability of this event by $P(B)$.
- (3) θ is incorrectly detected as being the destination of a flow originating in some predecessor of O . We denote the probability of this event by $P(C)$.

In cases 2 and 3, all the empirical probabilities are greater than or equal to τ , so

$$P_e \leq P(A) + P(B) + P(C) \tag{7}$$

In the following, we establish an upper bound for $P(B)$. The empirical rate on the path from O to j is given by $N_1(T)/T$.

$$P(B) \leq K e^{-\frac{N_1(T)}{T} T}$$

where K is a constant related to $\text{card}(G)$.

$$\frac{N_1(T)}{T} \geq r_k \geq p_k \tau \text{ for all } k \Rightarrow \frac{N_1(T)}{T} \geq \tau \frac{\sum p_k}{\# \text{ of links}} = \tau \tilde{p}.$$

Using the fact that $-\frac{N_1(T)}{T} \leq -\tau \tilde{p}$, we obtain

$$P(B) \leq K e^{-N_1(T)} \leq K e^{-\tau T \tilde{p}}$$

Similarly, for $P(C)$, we have $P(C) \leq K e^{-\tau T \tilde{p}}$.

For $P(A)$, the assumption is that there is at least one link along the path from O to θ with an empirical rate that is less than τ . Suppose j_0 is the first node on the path connected to a link whose rate is less than τ . Therefore, all links along the path up to j_0 have rates greater than or equal to τ . Let N_1 be the total number of timestamps taken during the observation period T .

We denote by λ the rate of the Poisson process defined by the recorded timestamps. Then, the rate on an arbitrary interval $[m_k, m_{k+1}]$ is $\lambda(m_{k+1} - m_k)$.

$$P(A) \leq \prod_{k=0}^{j_0-1} \frac{[r_k(m_{k+1} - m_k)]^{n_k}}{n_k!} e^{-r_k(m_{k+1} - m_k)} \prod_{k' \in N^+(j_0)} e^{-r_{k'}(m_{k'+1} - m_{k'})}$$

$N^+(j_0)$ denotes all successors of node j_0 on links with rates of at least R . Notice that in the first product $r_k \geq p_k \tau$ for all k , whereas all $r_{k'}$ from the second product are less than $p_{k'} \tau$. On the other hand, $r_{k'} \geq (N_1(T) - N(j_0)) / (T - m(j_0))$. We thus have

$$\begin{aligned} P(A) &\leq C_{j_0} \exp \left[-\tau \left(\sum_{k=0}^{j_0-1} p_k (m_{k+1} - m_k) \right) \right] \\ &\quad \exp \left[-\frac{N_1(T) - N(j_0)}{T - m(j_0)} (T - m(j_0)) \right] \\ &= C_{j_0} \exp(N(j_0)) \exp \left[-\tau \left(\sum_{k=0}^{j_0-1} p_k (m_{k+1} - m_k) \right) \right] \\ &\quad \exp \left(-N_1(T) \right) = K_1(j_0, \tau) \exp \left(-N_1(T) \right) \\ &= K_1(j_0, \tau) \exp(-T\rho(T)) \end{aligned}$$

where

$$C(j_0) = \prod_{k=0}^{j_0-1} \frac{r_k (m_{k+1} - m_k)^{n_k}}{n_k!}$$

and

$$K'(j_0, \tau) = C(j_0) e^{N(j_0) - \tau m_{j_0}}$$

Here, we let $\rho = \lim_{T \rightarrow \infty} N_1(T)/T > 0$.

Therefore, $P(A)$ decreases exponentially with the observation time.

Using Equation (7), we now have

$$\begin{aligned} P_e &\leq K'(j_0, \tau) e^{-\rho T} + \text{card}(G) e^{-\tau T \tilde{p}} + \text{card}(G) e^{-\tau T \tilde{p}} \\ &= e^{-T \tau \min\{\tilde{p}, \rho\}} (K'(j_0, \tau) + 2K) \end{aligned}$$

We can now investigate the convergence rate of our algorithm by showing the asymptotic decay rate of the error probability as time increases.

By logarithming the aforementioned inequality, we obtain

$$\lim_{T \rightarrow \infty} (-\ln P_e / T) \geq \tau \min(\tilde{p}, \rho) > 0$$

Compared with the upper bound for the error of algorithm TD, which is τ (corresponding to the case $p_k = 1$ for all k), the modified algorithm has upper bound $\tau \min(\tilde{p}, \rho)$.

The minimum is strictly positive as the minimum of a finite number of strictly positive values.

We can see that the algorithm has an exponential convergence rate on error probability over time T . That means, in our work, the error rate of finding correct connection pairs will decrease exponentially while the observing time goes up.

4.3. Generalization

We now prove a generalization of the previous result regarding the error probability when the number M of observable graphs goes to infinity.

If in the mobile case $\theta \neq \hat{\theta}$ for all M cases, then in each observable graph $G_i, i = 1, \dots, M$, we have $\theta \neq \hat{\theta}$. If we consider these M events independently, we have

$$P_e(T, M) \leq e^{-2TM} [2\text{card}(G)e^{\ln T + \ln 2} + K(j_0, 2)e^{-\tau T + \ln T + \ln 2}]^M = e^{-2TM} [2\text{card}(G)\tau T + K(j_0, 2)e^{(\tau-\rho)T}]^M$$

where $\rho(N, T) > 0$.

If we consider T to be fixed and sufficiently large and the number M of graphs going to infinity, the error probability again decays exponentially. We have

$$P_e(T, M) \leq (\text{card}(G)e^{-T \min(\rho, \tau\bar{p}, \tau\hat{p})})^M = (K'(j_0, \tau) + 2\text{card}(G))^M e^{-MT \min(\rho, \bar{p}, \hat{p})}$$

so P_e decreases with the speed of $e^{-\tau TM}$. We then have

$$\lim_{M \rightarrow \infty} [-\ln(P_e(T, M)/M)] \geq \lim_{M \rightarrow \infty} [\tau TM \min(\bar{p}, \hat{p}) + M \ln(K'(j_0, \tau) + 2\text{card}(G))] / M = \tau T \min(\bar{p}, \hat{p}) + \ln[K'(j_0, \tau) + 2K] > 0$$

We have obtained a bound on P_e , which is strictly smaller than 1 and thus nontrivial. It means that given a sufficient observation time, no matter how many the observation objects are, our tracing scheme still can obtain the correct connection pairs in a low error rate. In another words, it is scalable.

5. NUMERICAL EVALUATION

Given two connections C and C' , Figure 1 represents the frequency of occurrences of the same character a_i at every minute in the two connections over a time L . f denotes the

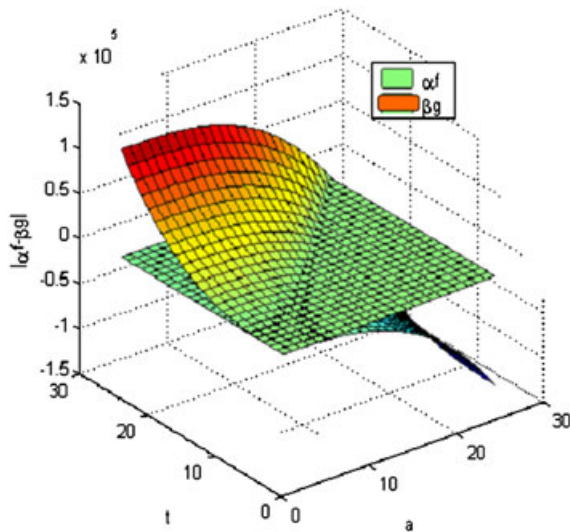


Figure 1. Maximum difference between f and g in absolute value.

frequency of a_i in connection C , and g represents the frequency of the same character in C' . The area between the graphs of f and g approximates the sum of the differences between the frequencies of characters in the two connections from Section 3: $B: \sum_{i=1}^L (f(a_i) - g(a_i))p(a_i)/L$.

Figure 2 represents the volume as given by the integral $\int\int_{[0, L] \times [0, T]} f(a, t) \alpha_k(a, t) da dt$. The volume of the solid in Figure 2 is bounded above by $f(a, t) \alpha_k(a, t)$ and below by $[0, L] \times [0, T]$ in the xy plane.

The volume of the solid in Figure 3 is bounded above by $f(Ly, T\tau) \alpha_k(Ly, T\tau)$ and below by $[0, 1] \times [0, 1]$ in the xy plane.

In Figure 3, the observation time L used for comparing two connections has been shrunk by a factor m . However, the number of samples (measured frequencies) taken is the same. Instead of measuring the frequency of a character every minute, we measure every $1/m$ minutes. The emphasized region determined by the graphs of f and g in Figure 3 has the same volume as the region in Figure 2. We can see how after the change of variable, the height of the figure increases as its base decreases, while the volume remains constant (Figures 2 and 3).

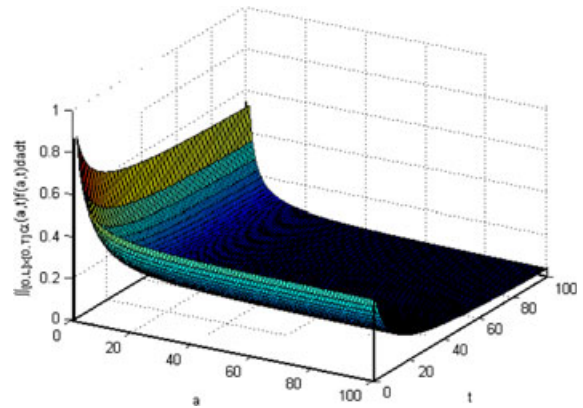


Figure 2. Volume of the solid that lies below function $f(a, t)\alpha_k(a, t)$ on $[0, L] \times [0, T]$.

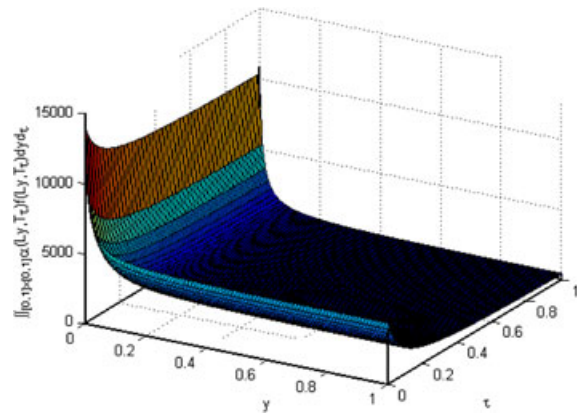


Figure 3. Volume of the solid that lies below function $f(Ly, T\tau)\alpha_k(Ly, T\tau)$ on $[0, 1] \times [0, 1]$.

An interesting consequence of the mean theorem ([42]) is shown in Figure 4. There exists a point $(x_0, y_0) \in S$ for which the volume of the integral $\iint_S f(x, y) dx dy$ is equal to $f(x_0, y_0) \text{Area}(S)$.

$f(x_0, y_0)$ gives the height of the rectangular parallelepiped whose volume approximates the value of the double integral $\iint_S f(x, y) dx dy$. This is particularly useful for finding the volumes of integrals that are hard to compute.

6. SIMULATIONS

To compare the performance of our approach to the original algorithms in [6], we adopt the same simulation method and settings in our experiment. Specifically, we deploy 25 nodes in an area of 750 m \times 750 m. Each node has a communication range of 250 m. All of the nodes move according to a slot-based random walk model. Within each slot (of length T seconds), every node will first randomly pick a point in the range of (0, 10] m on the basis of the previous position and then move toward it directly. Node 0 is always communicating with θ , and the other nodes randomly select a partner within each slot. We use the Dijkstra algorithm for path routing. Different paths may share one link and remain unchanged until the next slot. Every node generates traffic following a Poisson process at rate R . When the traffic passes by an intermediate node, it will be added with independent and identically distributed uniform delays drawn from $[0, \Delta]$ (the traffic order may be permuted because of queuing mechanism). After generating all the traffic, each link is padded with independent Poisson chaff noise to rate λ . We simulate the network for 1000 slots. In the simulation, $\lambda = 1$ packet/s, $\Delta = 1/18$ s, and $R = 1/9$ packet/s.

We have modified the TD and Recursive Tracing (RT) algorithms in [6] for our approach. The only difference between ours and the original is the selection of the support link. As we discussed in Section 4, assuming that the lack of information

(p_k) on each link is already known, the support link k is selected only if $r_k/p_k \geq \tau$. The modified algorithms are referred to as TD2 and RT2, accordingly. We simulate TD, TD2, and Simple Intersection Method (SIM) [6] under various thresholds and observation times. To find the best threshold τ , we first evaluate TD2 versus different τ values. As shown in Figure 5, for $M = 1$, we notice that the error probability of TD2 decreases with τ until it stabilizes around $\tau \in [0.19, 0.23]$ and then starts to increase sharply. By this observation, we set $\tau = 0.2$ in the subsequent tests to minimize error probability.

We then simulate TD2 and SIM [6] for various total observation times MT to evaluate their performance affected by the observation time. As we can see from Figure 6, just like the result of TD in [6], the error probability of TD2 also decays subexponentially. When M and T increase, the error probability of TD2 decreases. Because the SIM is not affected by the observation time, its error probability stays constant over range (0.55, 0.75).

In fact, as long as we have set up a good value for the threshold, the result of TD2 should have the same characteristics of that of TD. Hence, the only major difference between TD and TD2 is the selection of threshold. Finally,

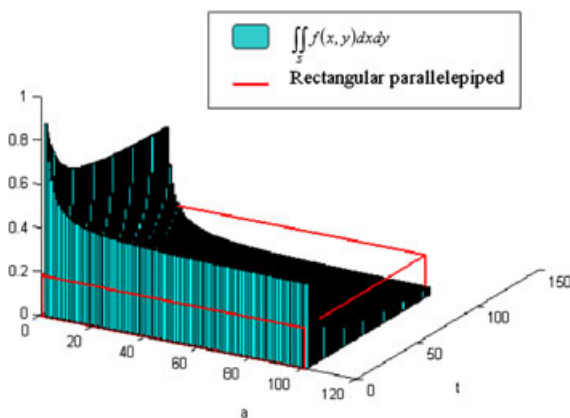


Figure 4. Approximating the volume defined by the double integral. The volume of $\iint_S f(x, y) dx dy$ is equal to the volume of the rectangular parallelepiped drawn in red.

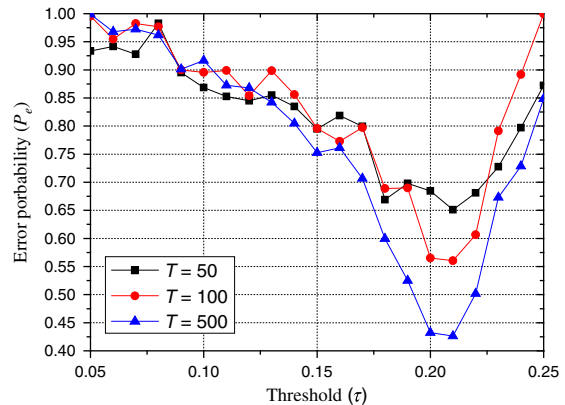


Figure 5. Error probability P_e of TD2 versus threshold τ .

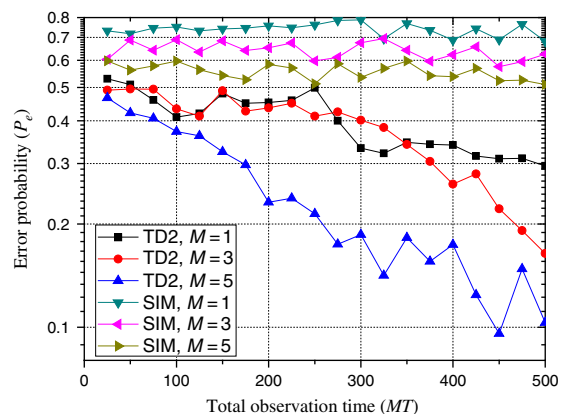


Figure 6. Error probability P_e of TD2 and SIM versus total observation time MT .

Table I. Comparison on performance between TD and TD2.

τ	$T = 50$		$T = 100$		$T = 500$	
	TD	TD2	TD	TD2	TD	TD2
0.12	0.5617	0.8451	0.5576	0.8539	0.4832	0.8679
0.13	0.5067	0.8551	0.4041	0.8987	0.4623	0.8422
0.14	0.6155	0.8348	0.4582	0.8559	0.4751	0.8048
0.19	0.6344	0.6979	0.5891	0.6899	0.5543	0.5252
0.20	0.6014	0.6844	0.5892	0.5652	0.5674	0.4325
0.21	0.6922	0.6511	0.5985	0.5604	0.5810	0.4263

we compare the performance of TD2 with the original TD in [6] under different threshold values. In Table I, we can see that a good threshold for TD ($\tau \in [0.12, 0.14]$) will lead to a high error probability (nearly 83%) for TD2, whereas a good threshold for TD2 ($\tau \in [0.19, 0.21]$) will lead a high error probability (nearly 60%) for TD. We therefore cannot tell which one is better on the basis of different thresholds that are predefined values. However, TD2 can achieve the lowest error probability when $\tau = 0.21$ and $T = 500$, which is about 4% lower than the TD's best result.

Note that the simulation procedure is an abstraction, varying from a simplified simulation to very detail emulation. The most detail one is an implementation. The chosen level of abstraction in simulation should be just enough to capture the essentials of the studied method. If the simulation is too detail (such as like an implementation), it not only costs much time but also obscure the effect of the studied methods because of the complexity of the system, which is treated as noise. If the simulation is too simple, it cannot study the method well. Therefore, a good approach is to choose an appropriate level of abstraction, neither too simple nor too detail, to study the method. We believe that our chosen simulation has already been good enough to evaluate the studied methods. To better illustrate the methods, the aforementioned deployed scenarios involve all the valuable information: random movement of each node, regular communication process, and background noise. In such context, we may reasonably draw the conclusion that our algorithm is an alternative way to trace targets in anonymous networks and that it is better than the algorithm in [6].

7. CONCLUSION

We provided a way to shorten the time interval used for thumbprinting and tune it suitably depending on network conditions. We have found interesting properties of the thumbprinting function by using the mean value and provided a general method to compute the value of the function. The method works for any function that satisfies certain properties. We have also studied another mechanism for tracing intruders in the Internet by passively monitoring flows between source and destination pairs. We computed the error probability of our algorithm and

showed that its value decreases exponentially as the observation time increases. Our simulation shows that our algorithm is an alternative way to trace targets in anonymous networks and that it is better than the algorithm in [6].

ACKNOWLEDGEMENT

This work is supported in part by the U.S. National Science Foundation, under grants CNS-0716211, CCF-0829827, CNS-0737325, and CNS-1059265. Prof. Chen's work is supported in part by the National Fundamental Research 973 Program of China under Grant 2011CB302801 and Macau Science and Technology Development Fund grant number 008/2010/A1.

REFERENCES

1. Internet, 2008, March 26. [Online]. Available: <http://en.wikipedia.org/wiki/Internet>
2. Staniford-Chen S, Heberlein LT. Holding intruders accountable on the Internet. In Proc. the 1995 IEEE Symposium on Security and Privacy, (Oakland, CA), May 1995, 39–49.
3. Jung HT, Kim HL, Seo YM, Choe G, Min SL, Kim CS. Caller Identification System in the Internet Environment. In Proc. of 4th UNIX Security Symposium, Oct. 1993, 69–78.
4. Wadell S. Private Communication, 1994.
5. Heberlein LT, Levitt K, Mukherjee B. Internet security monitor: an intrusion-detection system for large scale networks. In Proc. 15th National Computer Security Conference, Oct. 1992, 262–271.
6. He T, Wong HY, Lee K-W. Traffic analysis in anonymous MANETs. In Proc. IEEE MILCOM, San Diego, November, 2008.
7. Olteanu A, Xiao Y, Liu J, Chen TM. Studying non-intrusive tracing in the Internet. Proceedings of the 7th International ICST Conference on Heterogeneous Networking for Quality, Reliability, Security and Robustness, Qshine 2010.

8. Hamadeh I, Kesidis G. A taxonomy of internet traceback. *International Journal of Security and Networks* 2006; **1**(1/2):54–61.
9. Wang X. The loop fallacy and deterministic serialisation in tracing intrusion connections through stepping stones. *International Journal of Security and Networks* 2006; **1**(3/4):184–197.
10. Pan J, Cai L, Shen X. Vulnerabilities in distance-indexed IP traceback schemes. *International Journal of Security and Networks* 2007; **2**(1/2):81–94.
11. Korkmaz T, Gong C, Sarac K, Dykes SG. Single packet IP traceback in AS-level partial deployment scenario. *International Journal of Security and Networks* 2007; **2**(1/2):95–10.
12. Burt AL, Darschewski M, Ray I, Thurimella R, Wu H. Origins: an approach to trace fast spreading worms to their roots. *International Journal of Security and Networks* 2008; **3**(1):36–46.
13. Schrader KR, Mullins BE, Peterson GL, Mills RF. An FPGA-based system for tracking digital information transmitted via peer-to-peer protocols. *International Journal of Security and Networks* 2010; **5**(4):236–247.
14. Wong DS, Tian X. E-mail protocols with perfect forward secrecy. *International Journal of Security and Networks* 2012; **7**(1).
15. Vespa L, Chakrovorty R, Weng N. Lightweight testbed for evaluating worm containment systems. *International Journal of Security and Networks* 2012; **7**(1).
16. Kandah F, Singh Y, Zhang W, Wang T. A misleading active routing attack in mobile ad-hoc networks. *International Journal of Security and Networks* 2012; **7**(1).
17. Alsubhi K, Alhazmi Y, Bouabdallah N, Boutaba R. Security configuration management in intrusion detection and prevention systems. *International Journal of Security and Networks* 2012; **7**(1).
18. Xiao Z, Xiao Y. PeerReview re-evaluation for accountability in distributed systems or networks. *International Journal of Security and Networks* 2012; **7**(1).
19. Al-Salloum ZS. Defensive computer worms—an overview. *International Journal of Security and Networks* 2012; **7**(1).
20. Cao J, Jia X, Shu L. Editorial. *International Journal of Sensor Networks* 2012; **11**(1):1–2.
21. Silva R, Silva J, Caldeira JMLP, Rodrigues JJPC. Mobile multimedia in wireless sensor networks. *International Journal of Sensor Networks* 2012; **11**(1):3–9.
22. Medjiah S, Ahmed T, Hamid Asgari A. Streaming multimedia over WMSNs: an online multipath routing protocol. *International Journal of Sensor Networks* 2012; **11**(1):10–21.
23. Chia WC, Chew LW, Ang L, Seng KP. Low memory image stitching and compression for WMSN using strip-based processing. *International Journal of Sensor Networks* 2012; **11**(1):22–32.
24. Chew LW, Chia WC, Ang L, Seng KP. Low-memory video compression architecture using strip-based processing for implementation in wireless multimedia sensor networks. *International Journal of Sensor Networks* 2012; **11**(1):33–47.
25. Gamm GU, Kostic M, Sippel M, Reindl LM. Low-power sensor node with addressable wake-up on-demand capability. *International Journal of Sensor Networks* 2012; **11**(1):48–56.
26. Li D, Liu L, Du H. An approximation algorithm for dominating nodes selection in multi-channel multi-radio wireless sensor networks. *International Journal of Sensor Networks* 2012; **11**(1):57–65.
27. Takahashi D, Xiao Y, Meng K. Virtual flow-net for accountability and forensics of computer and network systems. *Wiley Journal of Security and Communication Networks*. DOI: 10.1002/sec.407, accepted.
28. Xiao Y, Yue S, Fu B, Ozdemir S. GlobalView: building global view with log files in a distributed/networked system for accountability. *Wiley Journal of Security and Communication Networks*. DOI: 10.1002/sec.374, accepted.
29. Xiao Y, Meng K, Takahashi D. Accountability using flow-net: design, implementation, and performance evaluation. *Wiley Journal of Security and Communication Networks*, Special Issue on Security and Privacy in Emerging Information Technologies 2012; **5**(1):29–49.
30. Zhang Y, Xiao Y, Ghaboosi K, Zhang J, Deng H. A survey of cyber crimes. *Wiley Journal of Security and Communication Networks* 2012; **5**(4):422–437.
31. Zhang Y, Xiao Y, Chen M, Zhang J, Deng H. A survey of security visualization for computer network logs. *Wiley Journal of Security and Communication Networks* 2012; **5**(4):404–421.
32. Xiao Y, Takahashi D, Liu J, Deng H, Zhang J. Wireless telemedicine and M-health: technologies, applications, and research issues. *International Journal of Sensor Networks (IJSNet)* 2011; **10**(4):202–236.
33. Liu L, Xiao Y, Zhang J, Faulkner A, Weber K. Hidden information in Microsoft word. *International Journal of Security and Networks (IJSN)* 2011; **6**(2/3):123–135.
34. Liu J, Xiao Y. Temporal accountability and anonymity in medical sensor networks. *ACM/Springer Mobile Networks and Applications (MONET)*, Special Issue: Wireless and Personal Communications 2011; **16**(6):695–712.
35. Takahashi D, Xiao Y, Zhang Y, Chatzimisios P, Chen H-H. IEEE 802.11 user fingerprinting and its applications. *Computers and Mathematics with Applications* 2010; **60**(2):307–318.

36. Liu J, Xiao Y, Ghaboosi K, Deng H, Zhang J. Botnet: classification, attacks, detection, tracing, and preventive measures. *EURASIP Journal on Wireless Communications and Networking* 2009, Article ID 692654:11. doi:10.1155/2009/692654.
37. Xiao Y. Flow-net methodology for accountability in wireless networks. *IEEE Network* 2009; **23**(5):30–37.
38. Meng K, Xiao Y, Vrbsky SV. Building a wireless capturing tool for WiFi. *Wiley Journal of Security and Communication Networks* 2009; **2**(6):654–668.
39. Lei M, Xiao Y, Vrbsky SV, Li C-C. Virtual password using random linear functions for on-line services, ATMs, and pervasive computing. *Computer Communications Journal*, Elsevier, Special Issue on Secure Multi-Mode Systems and their Applications for Pervasive Computing 2008; **31**(18):4367–4375.
40. Takahashi D, Xiao Y. Retrieving knowledge from auditing log files for computer and network forensics and accountability. *Wiley Journal Security and Communication Networks* 2008; **1**(2):147–160.
41. Xiao Y. Accountability for wireless LANs, ad hoc networks, and wireless mesh networks. *IEEE Communication Magazine*, Special Issue on Security in Mobile Ad Hoc and Sensor Networks 2008; **46**(4):116–126.
42. Mean value theorem, 2008, December. [Online]. Available: http://en.wikipedia.org/wiki/Mean_value_theorem
43. Implicit function theorem, 2007, October. [Online]. Available: http://en.wikipedia.org/wiki/Implicit_function_theorem
44. Poisson process, 2008, March 17. [Online]. Available: http://en.wikipedia.org/wiki/Poisson_process