

# Analyzing Privacy in Enterprise Packet Trace Anonymization\*

Bruno Ribeiro\*, Weifeng Chen<sup>+</sup>, Gerome Miklau\* and Don Towsley\*

\*Computer Science Department  
University of Massachusetts  
Amherst, MA, 01003

{ribeiro, miklau, towsley}@cs.umass.edu

<sup>+</sup>Computer Science Department  
John Jay College of Criminal Justice  
New York, NY, 10019

wechen@jjay.cuny.edu

UMass CMPSCI Technical Report 48-07  
September 2007

## Abstract

Accurate network measurement through trace collection is critical for advancing network design and for maintaining secure, reliable networks. Unfortunately, the release of network traces to analysts is highly constrained by privacy concerns. Several anonymization schemes have been proposed to address this issue. Preservation of prefix relationships among anonymized addresses is an important aspect of trace utility, but also causes a number of vulnerabilities in trace anonymization. In this work we present a novel, systematic attack on prefix-preserving anonymization which can be efficiently executed by an adversary in possession of a modest amount of public information about the network. The attack is general (encompassing a range of fingerprinting attacks proposed by others) and flexible (it can be adapted to emerging variants of prefix-preserving anonymization). Perhaps most importantly, we develop analysis tools that allow data publishers to quantify the worst-case vulnerability of their trace given assumptions about the adversary's external information. Using this analysis we quantify the trade-off between privacy and utility of alternatives to full prefix-preserving anonymization.

## 1 Introduction

Accurate network measurement through trace collection is critical for advancing network design and for maintaining secure, reliable networks. While the technological means exist to collect and analyze traces, the release of these traces to analysts is highly constrained by privacy concerns. Network trace data can include information about individuals (their personal data, communication habits, evidence of location), about an enterprise (its structure, organization, business practices), as well as about the underlying system (network topology, active services, security practices, etc.).

Because trace data is so sensitive, its publication and analysis is typically allowed only if the data is protected by an anonymizing transformation. Packet content is virtually always removed since its inclusion for unencrypted communication would constitute a severe privacy violation. The next step is to obscure source and destination IP addresses. This can be done with any one-to-one mapping, however the utility of trace data to researchers often requires that prefix relationships be maintained. For example, peer-to-peer systems measurement [1], worm outbreak analysis [2], the study of router forwarding caches, and prefix-based clustering on traces all depend on address locality and require the preservation of prefix relationships. Prefix-preserving anonymization was first implemented as a feature of the Unix tool `tcpdpriv` [3]. A number of other techniques have since been developed for efficient prefix-preserving packet trace anonymization [4–6].

---

\*This research has been supported in part by the NSF under grant awards CNS-0627642, ANI-0325868, and by CAPES (Brazil). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

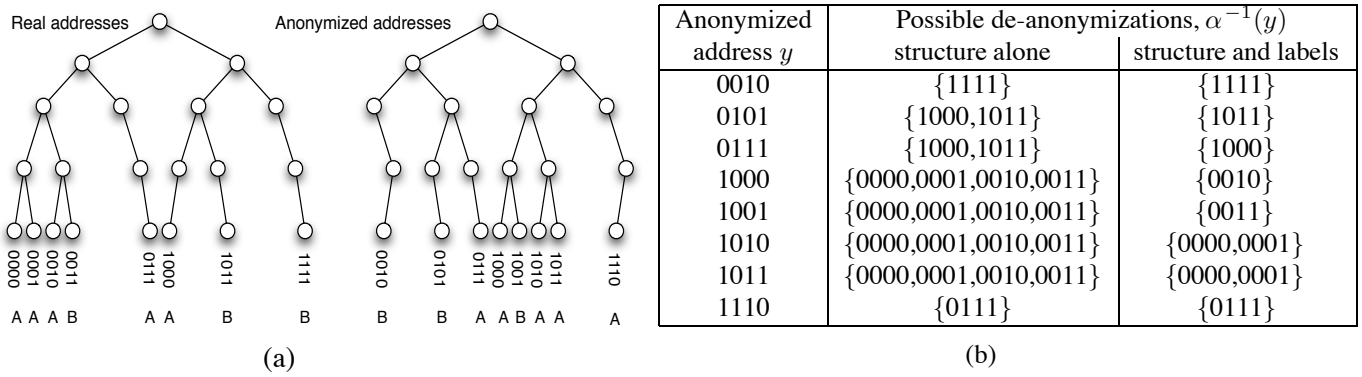


Figure 1: **(a)** A four-bit address space: (left) an address tree representing active hosts; (right) an anonymized address tree with prefixes preserved. **(b)** Anonymized addresses along with the set of possible de-anonymizations which can be inferred using structure alone (column 2) and using structure and attribute labels (column 3).

In this work we focus on enterprise IP traces; that is, traces collected near the perimeter of the Internet, at a gateway owned by an enterprise or institution. This is currently the most common type of trace available (because ISPs rarely release traces of backbone traffic). Numerous enterprise IP traces have been released in anonymized form [7–13].

Our work presents a novel, systematic attack on prefix preserving anonymization which can be efficiently executed by an adversary in possession of a modest amount of public information about the network. The attack has a number of advantages over others recently proposed. In particular, we provide an analysis tool that allows data publishers to quantify the worst-case vulnerability of their trace given assumptions about the adversary’s external information. We demonstrate the effectiveness of the attack and analysis tools on a real enterprise trace.

## Background and Related Work

Active addresses in an enterprise network with  $2^k$  total IP addresses can be represented as leaves in a binary tree of height  $k$ , which we call an *address tree*. Figure 1(a) (left) shows an example of an address tree in a 4-bit address space. An anonymization function  $\alpha$  is a bijective function taking an IP address  $x$  to its anonymized counterpart  $\alpha(x)$  anywhere in the IPv4 address space. An anonymization function *preserves prefixes* when two addresses  $x$  and  $z$  agree on the first  $i$  bits if and only if  $\alpha(x)$  and  $\alpha(z)$  also agree on their first  $i$  bits. Figure 1(a)(right) shows one possible prefix-preserving anonymization of the original address tree.

**Attack overview** Our attack focuses on the central threat in trace publication: *host de-anonymization*, in which an adversary is able to associate the actual IP address,  $x$ , with the obscured IP address,  $\alpha(x)$ , present in a published trace. The attack is based on the observation that if the adversary can gather information about the existence of addresses in the real network, a set of structural constraints emerge that lead to serious disclosures. This attack is most effective when applied to the internal addresses of an enterprise trace, which are easily distinguishable from external addresses, and whose true identity are often the most sensitive to the enterprise.

For example, assume the adversary is able to discover all active addresses in the real network, as illustrated in Figure 1(a)(left). Using structural relationships alone, the adversary can de-anonymize some of the addresses shown in the anonymized tree. For each anonymized address, its possible de-anonymizations are shown in Figure 1(b). Two anonymized addresses (0010 and 1110) are uniquely de-anonymized based on the structure of the network alone. Intuitively, address 1110 is distinguishable because it is alone in its 2 bit subtree, but it is adjacent to a complete subtree of 4 nodes. For other addresses, unique de-anonymization is not possible and a set of possible candidates remain, as shown in Figure 1(b). For each anonymized address  $y$  in the  $\{1000, 1001, 1010, 1011\}$ , its de-anonymization  $\alpha^{-1}(y)$  must be one of  $\{0000, 0001, 0010, 0011\}$ . These addresses remain hidden in a crowd.

Continuing the example, if the adversary has additional information about properties of real addresses, the possible de-anonymizations are further constrained. For example, the adversary may know that host 0011 has port 22 open and witnesses frequent traffic on this port. Traffic on port 22 can be observed in the trace for some hosts but not others. The host properties (known about real hosts, observed for trace hosts) can be represented as labels on the leaf

nodes (in Figure 1(a) the labels are simply A and B). Assuming only addresses with common labels should match the possible de-anonymizations are further limited, as shown in the final column of Figure 1(b). Address 1001 and its sibling 1000, previously hidden in a crowd of size 4, are now uniquely de-anonymized.

The effectiveness of this attack depends, first, on the completeness and accuracy of an adversary’s external information about hosts in the real network (whether they are active, properties of the traffic, etc.) and second, on properties of the network itself (the allocation of addresses, the uniformity of host properties). As suggested by this example, the adversary’s inference process involves matching the leaves of the address trees in a way that respects constraints imposed by the tree structure as well as the host labels. In practice, an adversary’s information will be incorrect for some hosts. As a result, the attack algorithm we describe is based on cost-based approximate tree matching [14, 15].

**Related Work** Vulnerabilities in IP trace anonymization are by now widely recognized. Slagell et al. [16] provide a thorough categorization including *injection attacks* (in which a recognizable pattern of traffic is introduced to the trace by the attacker), *fingerprinting attacks* (in which properties of real addresses are matched to addresses exhibiting those properties in the trace), and *structure recognition* (in which, for example, the de-anonymization of one address is used to narrow the possible de-anonymizations of other addresses when prefix structure is preserved).

Xu et al. [5] evaluated the impact of structure recognition attacks from small sets of de-anonymized hosts by calculating the overall average number of discovered bits in the anonymization. Koukis *et. al.* [17] proposed two attacks on anonymized traces (although they do not capitalize on prefix preservation). The first is a fingerprinting attack that matches profiles of web site content with traffic present in the trace in order to identify well-known websites. (The hosts in the trace visiting these sites remain protected.) The second assumes the attacker can identify existing patterns of linear port-scanning activity. Brekne et al. [18] describe both fingerprinting and injection attacks on prefix-preservation for powerful adversaries assumed capable of forging trace traffic and granted knowledge of the traffic distribution.

Despite these vulnerabilities, there are few techniques designed to resist them. One notable exception is the work of Pang et al. [19], in which scans are removed from traces to mitigate injection attacks, and prefixes are not preserved for addresses internal to the enterprise. The choice to forgo full prefix-preservation for internal addresses sacrifices the utility of the trace in order to resist attacks like the one presented in this work and others mentioned above. Instead of pure prefix preservation, the subnet and host portions of each address are anonymized separately, but subnet relationships are preserved. We refer to this as an example of *partial prefix preservation*.

Recently, Coull *et. al.* discussed attacks on this presumably safer anonymization. They describe a fingerprinting attack which creates a behavioral profile for hosts using public information sources (e.g., DNS or web search engine queries), the perceived popularity of the target host, and its possible locations within the network topology. Some parts of the network topology hidden by the anonymization are recovered and some distinguished public servers are identified.

**Contributions** The existing attacks on prefix-preserving anonymization suffer from two shortcomings which our techniques address. First, most fingerprinting attacks attempt re-identification on each host independently, without reasoning collectively over the entire trace. Second, the effectiveness of the attacks is highly dependent on adversary knowledge which is nearly impossible for trace publishers to estimate. We make the following main contributions to address these issues:

- We describe our attack (Section 2), which systematically combines fingerprinting and structure recognition attacks. It can easily accommodate any form of external information (by representing it as labels on address nodes) including information acquired through injection attacks. It can be applied to both full or partial prefix preservation, and the attack is efficient: it takes a few minutes on common hardware to execute the attack on a class B network.
- In Section 3 we perform a thorough experimental evaluation of the attack on a trace collected at a large university. We show that an adversary, using probes of 9 distinct TCP port numbers, can *uniquely* re-identify 17% of the active hosts in the trace, while about 50% of the hosts have candidate sets of size less than 8.

---

<sup>1</sup>This assumption does not always hold. We consider this issue fully in Section 3.

- We analyze our attack formally in Section 4 by proving that the worst case rate of de-anonymization (for a given set of observable host properties) is achieved by an adversary whose external information is perfectly accurate.
- In Section 5 we evaluate the impact of the partial prefix preservation recommended in [19]. Our techniques allow us to quantify the improvement in privacy gained through this technique, but we also show that one aspect of it (randomization of subnets) sacrifices trace utility without increasing privacy.

The worst-case analysis of our generic fingerprinting attack has important practical consequences. It can form the basis for an efficient tool that data publishers can use to conservatively assess the risk of publishing traces for their particular network, to identify the most vulnerable hosts, or even to guide address allocation within their network.

## 2 De-anonymization attack

In this section we describe our attack on prefix-preserving anonymization and its connection with approximate tree matching algorithms. We assume in this section that the attack is applied to full prefix preservation and return to partial prefix preservation in Section 5. The attack consists of three major steps:

**Step I** The adversary derives traffic fingerprints for each host in the anonymized trace.

**Step II** The adversary collects information from external public sources to construct fingerprints for network hosts.

**Step III** The adversary uses the fingerprints obtained in steps I and II above to recover the anonymization function,  $\alpha$ , in partially or in full.

### 2.1 Attack description

**Step I: Deriving trace fingerprints.** There are many types of traffic information that can be recovered from an anonymized packet trace. For instance, an address  $y$  is “active” when the trace contains at least one packet with source address  $y$ . Adversaries can also gather information about which TCP services address  $y$  provides by examining the trace for packets with source address  $y$  from well known service port numbers. Information on host operating system types can also be obtained [20]. Other types of information include flow sizes, packet sizes, and packet inter-departure times. Appendix A includes more details on such attributes. We refer to the collection of traffic information about an anonymized address  $y$  as the *trace fingerprint* of  $y$ .

**Step II: Mining external fingerprints.** In this second step, the adversary gathers external information about the real network of hosts believed to be present in the trace. The host properties collected are similar to the those presented in Step I. The adversary may gather information from reverse DNS queries, URL analysis, Web-crawling, active probing, or other public sources. Network service port status, such as active/inactive Web and E-mail TCP ports, are among the easiest types of external information to gather from a network. Such information can be obtained by probing the network with TCP SYN packets sent to ports such as 25 (E-mail) or 80 (Web). The adversary infers that IP address  $x$  runs a network service at the Web server port when  $x$  replies with a TCP SYN ACK when probed on port 80.

In what follows we provide a formal definition for the above fingerprints. Let  $\Gamma_i$  be an alphabet of symbols (or labels). Let  $Y_i(y) \in \Gamma_i$  be a variable that refers to a traffic characteristic of anonymized address  $y$ . For instance,  $\Gamma_i = \{\mathbf{A}, \mathbf{I}\}$  and  $Y_i(y) = \mathbf{A}$  when  $y$  is “active” or  $Y_i(y) = \mathbf{I}$  when  $y$  is “inactive”. Let  $F_t(y)$ , the *trace fingerprint* of  $y$ , be a vector with  $k$  of these variables:  $F_t(y) = (Y_1(y), \dots, Y_k(y))$ . Let  $\mathcal{A}$  be the set of all anonymized addresses and  $F_t = \{(y, F_t(y)) : \forall y \in \mathcal{A}\}$  be the set of all fingerprints of a trace. Let  $x$  be an un-anonymized address of the network. External fingerprints  $F_e(x) = (X_1(x), \dots, X_k(x))$  are constructed in the same way as  $F_t(y)$ . Throughout this work we assume that  $X_i$  and  $Y_i$  refer to the same network characteristic such as “active address” or “TCP Web service”.

*Inaccuracy of external fingerprints.* In real world scenarios, adversaries are unlikely to collect perfectly accurate external fingerprints. We say that an adversary is able to obtain *perfect external fingerprints* when for all  $x$ ,  $F_e(x) = F_t(\alpha(x))$ ; otherwise, fingerprints are *imperfect*. Firewalls and changes to the network are among the most common reasons for these mismatches. Firewalls can prevent adversaries from collecting reliable address attributes using active probing. Because we do not assume adversary probing is synchronized with trace collection, changes to the network

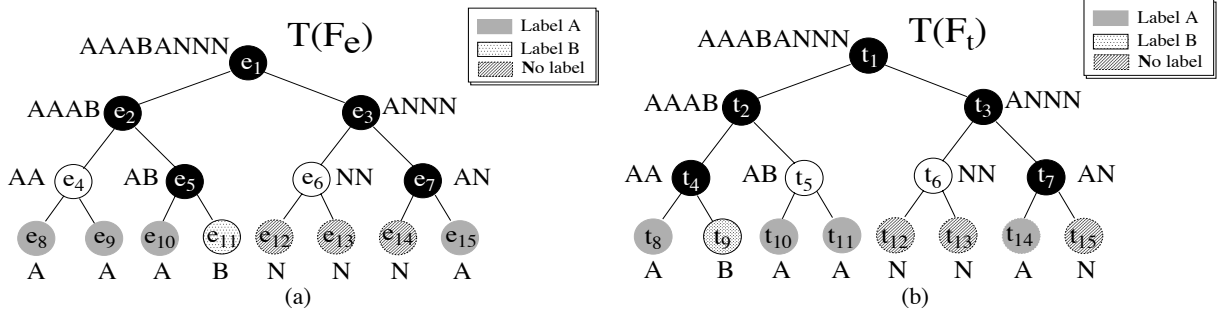


Figure 2: **(a)** External fingerprint tree  $T(F_e)$  obtained from Figure 1(a)(left) example; **(b)** Trace fingerprint tree  $T(F_t)$  obtained from Figure 1(a)(right) example. Symbols A and B represent the labels of Figure 1(a) and symbol N represents an address with no label.

can occur. Appendix B presents a compilation of fingerprint inconsistencies, each of them with a possible countermeasure that can be taken by an adversary. The compilation in Appendix B is far from complete but gives an idea of the challenges that an adversary faces.

Further, imperfect external fingerprints require that we perform an approximate matching between trace and external fingerprints. As a consequence, it is necessary to define a cost function which is applied to pairs of fingerprint attributes and acts as a measure of the adversary’s certainty about the association between a trace and an external fingerprint.

**Definition 2.1 (Cost function).** Let  $c(F_e(x), F_t(y)) \geq 0$  be a cost function that is zero when  $F_e(x) = F_t(y)$ .

**Step III: Recovering the anonymization function.** The basic idea behind our attack is to find a set of de-anonymization functions that are “good approximations” to the correct de-anonymization function  $\alpha^{-1}$ . In what follows we formally define the above as the adversary’s objective.

**Definition 2.2 (Adversary’s objective).** Let  $\mathbb{A}$  be the set of all valid de-anonymization functions. Let

$$\gamma(\tau, F_e, F_t, c) = \sum_{\forall y \in \mathcal{A}} c(F_e(\tau(y)), F_t(y)) \quad (1)$$

be the cost of matching fingerprint  $F_t(y)$  to its de-anonymized counterpart  $F_e(\tau(y))$  using the de-anonymization function  $\tau \in \mathbb{A}$ . The objective of the adversary is to find the set of all de-anonymization functions  $\mathbb{B}_{F_e, F_t, c} \subseteq \mathbb{A}$  that have the smallest cost  $\gamma$  for fingerprints  $F_e$  and  $F_t$ ; or

$$\mathbb{B}_{F_e, F_t, c} = \operatorname{argmin}_{\tau \in \mathbb{A}} \gamma(\tau, F_e, F_t, c). \quad (2)$$

Note that the success of the adversary (in finding a good approximation to  $\alpha^{-1}$ ) depends on both the fingerprints obtained from Steps I and II ( $F_t$  and  $F_e$  respectively), and  $c$  (Definition 2.1). We can measure the adversary’s success in attacking an anonymized address  $y$  using fingerprints  $F_e, F_t$  and function  $c$  using

$$\beta(y, F_e, F_t, c) = \{\tau(y) : \forall \tau \in \mathbb{B}_{F_e, F_t, c}\}, \quad (3)$$

which is the set of good matches of  $y$  (or the *match set* of  $y$ ). In what follows we show an algorithm that computes  $\mathbb{B}_{F_e, F_t, c}$  (equation (2)). Note that  $\mathbb{B}_{F_e, F_t, c} \subseteq \mathbb{A}$  and therefore  $\mathbb{B}_{F_e, F_t, c}$  contains only valid de-anonymization functions. Thus, Step III must enforce matchings to be consistent with the (full or partial) prefix-preserving order. For this we introduce the use of *fingerprint trees*.

**Definition 2.3 (Fingerprint tree).** Let  $F$  be a set of fingerprints. A fingerprint tree  $T(F)$  of a set of fingerprints  $F$  is an address-space tree in which each node is assigned a label. The label of the leaf corresponding to address  $y$  is a string of symbols  $f(y) = Y_1(y) \dots Y_k(y)$ . Inner vertex labels are defined by the following recursion: Let  $y$  and  $z$  be two sibling vertices in the fingerprint tree and  $a$  be their parent. Vertex  $a$  is labeled

$$f(a) = \begin{cases} f(y) \cdot f(z) & \text{if } f(y) <_{LEX} f(z), \\ f(z) \cdot f(y) & \text{otherwise,} \end{cases}$$

where “.” represents a string concatenation and  $a <_{LEX} b$  if  $a$  is less than  $b$  in lexicographical order. We call inner vertex  $a$  “white” if  $f(y) = f(z)$ , otherwise it is called “black”.

There are two types of fingerprint trees: *trace fingerprint trees*,  $T(F_t)$ , and *external fingerprint trees*,  $T(F_e)$ . We exemplify these two fingerprint trees by constructing them from the labels given in Figure 1(a). Figure 2(a) shows the subtree of  $T(F_e)$  that encompasses addresses 0000 to 0111 of Figure 1(a)(left). Figure 2(b) shows their corresponding trace fingerprints. Inner vertices of these fingerprint trees are painted white and black according to the vertex denomination given in Definition 2.3.

## 2.2 Attack algorithm

Let  $T(F_t)$  and  $T(F_e)$  denote trace and external fingerprint trees created using fingerprints  $F_t$  and  $F_e$  respectively.  $T(F_t)$ ,  $T(F_e)$ , and a cost function  $c$  are given as inputs to Step III, which outputs  $\beta(y, F_e, F_t, c)$  for all anonymized addresses in the trace. The attack is quite straightforward. The above problem can be formulated as a trivial extension of a constrained tree edit distance problem for unordered trees, described as follows. Consider relabeling tree  $T(F_e)$  such as to make the complete binary trees  $T(F_t)$  and  $T(F_e)$  isomorphic. Now consider that each leaf relabeling operation has a non-negative cost associated with it and relabeling of inner vertices have cost zero. Tree edit distance algorithms find all sets of relabeling operations over  $T(F_e)$  with minimum total cost as to make  $T(F_t)$  and  $T(F_e)$  isomorphic. Using the above formulation, our problem is an instance of the tree alignment distance problem [15]. Our problem also has more specific constraints that can be used to optimize the minimum cost search. These other constraints are:  $T(F_t)$  and  $T(F_e)$  are complete binary trees and edit operations are restricted to relabeling operations (no insertion or removal of vertices). Using these restrictions we provide a fairly efficient algorithm for instances of the problem that the adversary is likely to encounter. Appendix C carries an example of how our algorithm works. The Java source code of our implementation of the above algorithm is also available for download<sup>2</sup>.

**A remark on the choice of cost function** A cost function that works well despite the level of external fingerprint imperfectness is likely to require knowledge of  $\alpha$  beforehand, which is exactly what our attack is looking for. However, it is possible to learn better editing costs during the search for the edit distance matches [21]. The algorithm presented in [21] is exponential on the size of the network and thus unfit for our purposes. Further research is need to determine whether learning editing costs can be made practical for our setting. However, we present experimental evidence suggesting that the choice of cost function is not too critical for an effective attack.

In the following section we apply our attack against real anonymized traces of a large university network. The following results were obtained for a network with 65536 addresses. In our experiments the algorithm typically finishes within a couple of minutes on a 1.83GHz Intel Centrino Duo processor.

## 3 Experimental results

The attack presented in Section 2 has three parameters:  $F_t$ , the trace fingerprints,  $F_e$ , the external fingerprints, and  $c$ , the cost function. This section explores the effectiveness of our attack and the importance of these parameters for real prefix-preserved anonymized traces. The following experiments are performed over full prefix-preserved traces collected at an university Internet gateway. The observed network has 65536 addresses. The external fingerprints,  $F_e$ , are gathered through network probing from an external host. Note that our fingerprint attack is not an injection attack. Our attack is passive, i.e., we do not assume that the anonymized trace contains any information about our probes. Throughout this section we focus solely on the de-anonymization of addresses that are internal to the network.

**Trace fingerprints** Our traces were collected at an Internet gateway of a class B university network. The university has two other gateways connected to the Internet that were not monitored. We attack five 24 hour traces collected between June 18th to June 22nd, 2007 from 12:00AM to 11:59PM each. In this section we present the representative results of the June 18th trace, named here *Trace-0618*. This trace has 573,037,780 packets, 9097 active internal network addresses, and was anonymized using prefix-preserving anonymization before its release.

---

<sup>2</sup><http://www-net.cs.umass.edu/~ribeiro/deanonymization/>

Following **Step I** of our attack, we derive the following fingerprint attributes for each address in the trace. The “Active” attribute is set “true” for an address  $y$  if the trace has at least one packet with source address  $y$ . Nine more attributes are derived from the existence of at least one TCP SYN ACK packet for address  $y$  on ports corresponding to common services: FTP, SSH, Telnet, E-mail, Time, DNS, Web, POP3, and SOCKS. A final attribute (TTL) is derived from the time-to-live IP field of traffic for address  $y$ . Popular host operating systems have distinct initial TTL values, and we distinguish between four main TTL initial values: Windows, Linux, MacOS, and “Other”. Any inconclusive labels are assigned “undefined”. Thus we consider 11 fingerprint attributes in total. Others are of course possible.

**External fingerprints** Following **Step II** of our attack, we actively probe all addresses in network from an external host (in Brazil). We do not assume the adversary knows when traces are being taken, so in general there will be a temporal mismatch between the time of trace collection and the time of adversary probing. In order to test the impact of time in the accuracy of the de-anonymization, we collect fingerprints on June 14th (*External-0614*), June 15th (*External-0615*), June 18th (*External-0618*), July 19th (*External-0719*), and August 27th (*External-0827*) of 2007. We were careful to remove the probes of *External-0618* from the trace *Trace-0618* in order to avoid introducing a bias to our measure. External fingerprint attributes correspond directly to trace fingerprint attributes and represent the network characteristics “Active”, FTP, SSH, Telnet, E-mail, Time, DNS, Web, POP3, SOCKS, and initial TTL value. The attribute “Active” represents the absence or presence of any response to the probes from a particular address.

**Cost function** Unless stated otherwise, we use a cost function where the cost of editing  $Y_i(y)$  to match  $X_i(\tau(y))$  is not independent of the value of  $Y_i(y)$ . We define this type of cost function as an *asymmetric cost* function. Our cost function reflects the natural belief that, for example, it is more likely to find an open port 80 when probing  $\alpha^{-1}(y)$  without recorded traffic on port 80 in the trace, than finding attribute traffic on port 80 in the trace but port 80 is closed during the probing. The value of each editing operation is 1, except in the case where traffic is observed in the trace and the corresponding TCP port was not open. In this case the editing cost is zero. The cost of editing an “undefined” label is also zero. Our experiments did not indicate significant sensitivity to symmetric or asymmetric cost functions. For lack of space we omit the experiments demonstrating stability under these variations of the cost function.

### 3.1 Result analysis

To assess the overall quality of the de-anonymization we use two metrics. The first metric is the set  $M(K, F_e, F_t, c)$ , which measures which addresses  $y$  have match sets  $(\beta(y, F_e, F_t, c))$  of size no greater than  $K$ . More formally let  $\mathcal{A}_a$  be the set of anonymized addresses in the trace that are active and

$$M(K, F_e, F_t, c) = \{y : \forall y \in \mathcal{A}_a; \text{ and } |\beta(y, F_e, F_t, c)| \leq K\}.$$

As  $F_e$  can have imperfect information, some of the addresses in  $M$  can have match sets that do not contain the correct de-anonymized address, i.e.,  $\alpha^{-1}(y) \notin \beta(y, F_e, F_t, c)$ . Thus, we would also like to know which addresses in  $M$  have match sets that contain the correct de-anonymized address. For this we introduce a second metric

$$V(K, F_e, F_t, c) = \{y : \forall y \in M(K, F_e, F_t, c); \text{ and } \alpha^{-1}(y) \in \beta(y, F_e, F_t, c)\}. \quad (4)$$

We also refer to  $V(K, F_e, F_t, c)$  as the set of  $K$ -vulnerable anonymized addresses with respect to parameters  $F_e$ ,  $F_t$ , and  $c$ . An anonymized address  $y$  is  $K$ -vulnerable with respect to  $F_e$ ,  $F_t$ , and  $c$  if  $|\beta(y, F_e, F_t, c)| \leq K$  and  $\alpha^{-1}(y) \in \beta(y, F_e, F_t, c)$ . Note that  $|M(K) - V(K)|$  shows the number of “errors” our attack makes. Also note that  $M$  and  $V$  are cumulative, i.e.,  $V(K) \subseteq V(K + 1)$ . Ideally we would like to have  $V(K)$  large with  $M(K) - V(K)$  small for any value of  $K$ . In what follows we omit the dependency of  $V$  and  $M$  on  $F_t$ ,  $F_e$ , and  $c$ . The choices of  $F_t$ ,  $F_e$ , and  $c$  should be clear from the context.

Following **Step III**, our first experiment uses external fingerprints obtained from *External-0618* and trace fingerprints obtained from *Trace-0618*. Figure 3(a) shows the values of  $|M(K)|$  and  $|V(K)|$  of this experiment. We observed that  $V(1) = 1620$ , i.e., 1620 active anonymized addresses, 17% of the number of active addresses in the trace, had their real addresses disclosed. The number of errors is  $|M(1) - V(1)| = 264$  from which we conclude that our attack is fairly accurate. In this experiment, all active addresses had match sets of size at most 1024. A little

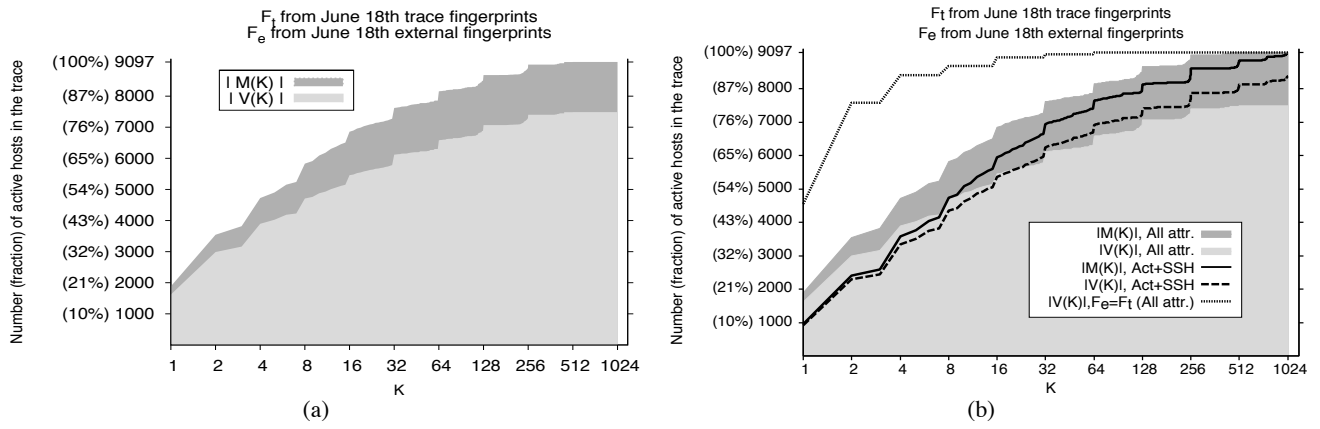


Figure 3: **(a)** The  $x$ -axis of our graph shows  $K$ , the size of the match sets. The lighter area shows  $|V(K)|$  and the darker area shows  $|M(K)|$ .  $|M(K) - V(K)|$  is the number of de-anonymization errors; **(b)** This figure shows five curves. The two solid areas show the curves of Figure 3(a). The curve “ $|V(k)|, F_e = F_t(\text{All attr.})$ ” considers the case where external fingerprints are perfect. The remaining curves consider fingerprints with attributes “active” and “SSH”.

more than 50% of the active addresses had matches of size no greater than 4. Among these matches, 85% of them are correct, i.e.,  $\alpha^{-1}(y) \in \beta(y, F_e, F_t, c)$ .

Note that although the percentage of hosts uniquely identified is a small portion of the total trace nodes, this represents a very significant vulnerability. Trace anonymization is intended to conceal the true identity of all nodes present in the trace. Instead, a significant portion of them are uniquely identified, and it is possible for the adversary to identify a small set of candidate matches for a much larger percentage of nodes. If an attacker wishes to re-identify a specific trace host, reducing the set of feasible candidates to 8 or even 16 may be more than sufficient. The attacker could then acquire more specific information to refine their fingerprints and also eliminate any false re-identifications that may be present. Overall, these results show that prefix-preserving anonymization is not safe against this type of attack. In the next two subsections we analyze the relative impact of different fingerprint attributes and the impact of collecting external fingerprints before or after trace collection.

### 3.2 Importance of distinct fingerprint attributes

In what follows we look at the sensitivity of our attack against a number of external fingerprint attributes from *External-0618*. We extract some traffic attributes from *Trace-0618* and *External-0618* and present the attack results using these attributes as fingerprints at Figure 3(b). Figure 3(b) shows five curves. The two solid areas show the curves of Figure 3(a). The curve “ $|M(K)| \text{ Act+SSH}$ ” shows results of the same attack now only using fingerprints attributes Active and SSH. Note that removing all other attributes from the fingerprints tend to increase the size of the matches. However, from  $|V(K)|$  we see that the fraction of matches that are correct is much higher in this scenario (95% of the matches of  $M(1)$  are correct, against 85% obtained with all attributes). Further analysis of our data reveals that adding the TTL type attribute to the fingerprint significantly increases  $M(K)$  but also increases the fraction of errors for small values of  $K$ . All other attributes cause minor impacts on  $M$  and  $V$ . From our analysis we conclude that removing the TTL attribute improves the overall accuracy of our attack. Here we also compare our results to the case where external fingerprints are perfect, i.e.,  $F_e = F_t$ , using attributes “Active+SSH”. The curve “ $|V(k)|, F_e = F_t(\text{A+SSH})$ ” shows the case where external fingerprints are perfect. We study this case in detail in Section 4. Note that in this case  $M(K) = V(K)$ . We can see that the noise introduced by imperfect external fingerprints significantly reduces the accuracy of our attack.

### 3.3 Attack sensitivity to late or early probing

Because the adversary does not know the date and time of trace collection, their fingerprint gathering will not necessarily occur near that time. The passage of time can affect accuracy because of changes to hosts, services running on hosts, or to address allocation policies. We find in our experiments that probing for external fingerprints as long



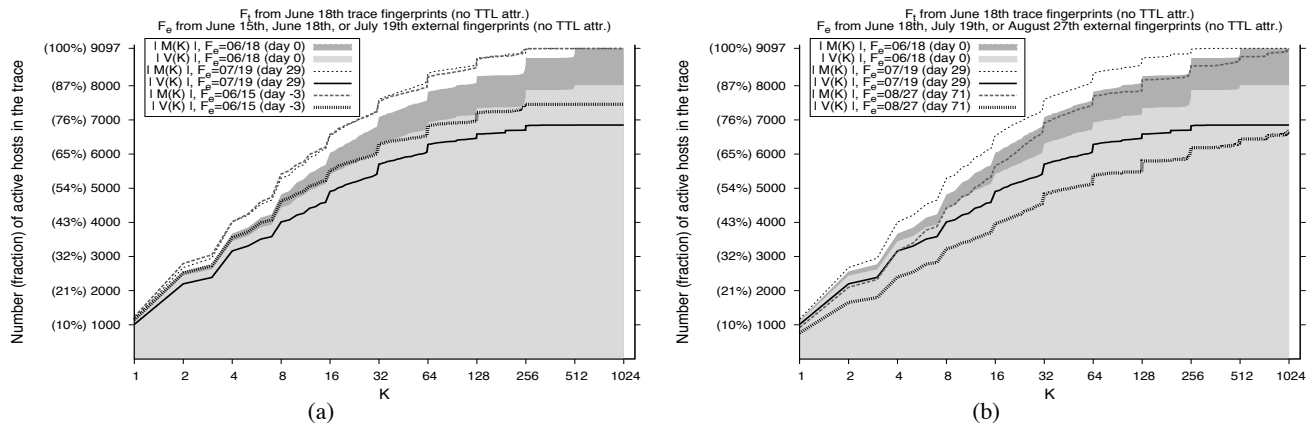


Figure 4: (a) These six curves show the values of  $|M(K)|$  and  $|V(K)|$  for all attributes except “TTL”. The two solid lines represent external fingerprints collected at the day of the trace collection (“day 0”). External fingerprints of the remaining two pairs of curves were collected three days before the trace collection (“day -3”) and 29 days after the trace collection (“day 29”) respectively; (b) This figure complements Figure 4 (a) with external fingerprints collected 71 days after (“day 71”) the trace collection.

as a month after trace collection leads to reasonable results but that accuracy does noticeably decline after about two months.

In the following experiments we use external fingerprints (all fingerprints except the initial TTL attribute) from *External-0615*, *External-0618*, *External-0719*, and *External-0827*. We drop the “TTL” attribute and use only the most perfect external attributes in order to study the impact of the passage of time on our results. Figures 4(a) and 4(b) plot  $|V(K)|$  against  $K$  for the external fingerprints obtained at these dates. We start comparing the matching results from *External-0618* (day 0) and *External-0615* (day -3). Probing 3 days before the trace was collected increases the number of matching errors (match sets without the correct de-anonymized address) but does not reduce the number of correct matches. However, considering only small match sets (e.g.  $K = 1$ ), the absolute number of matching errors is still small. Thus the attack is still fairly accurate. Comparing the matching results from *External-0618* (day 0) and *External-0719* (day 29) we see that probing 29 days after the trace is collected significantly increases the number of matching errors. It also slightly reduces the number of correct matches. Once again the absolute number of matching errors is still fairly small for small values of  $K$  and  $|M(K)|$  is still high. Again, the attack is still fairly efficient. Comparing *External-0618* and *External-0827* (Figure 4(b)) shows that matching errors are four or more times greater when external fingerprints are probed 71 days after the trace is collected. The absolute number of small match sizes also reduces. The fraction of matches that are incorrect with  $K = 1$  is still small, a little less than 20%, and  $|M(1)| = 917$  is still large. These results suggest that the adversary has a fairly large time window to collect useful external fingerprints.

## 4 Worst case analysis of host de-anonymization

The experimental results presented in the previous section show that the effectiveness of the attack depends on the accuracy of the external fingerprints collected by the adversary. For hosts that remain hidden, it is not clear whether their anonymity is due to the trace transformation techniques, or due to the weakness on the part of the adversary. In this section we analyze our attack under the simple but important assumption that the adversary’s external fingerprints are *perfect*, that is, they match precisely the properties of the trace nodes. For any fixed set of fingerprint attributes under consideration, we show formally that no adversary could be *more successful* in host de-anonymization (independent of the choice of cost function). This allows data publishers to efficiently calculate, prior to publication, the worst-case disclosure for a given trace and to pinpoint vulnerable nodes.

We apply this analysis to the experimental trace from the last section showing the unique de-anonymizations jump from 17% to 44% of trace nodes. In addition, we apply the worst-case analysis to partial prefix preservation, with some surprising consequences.

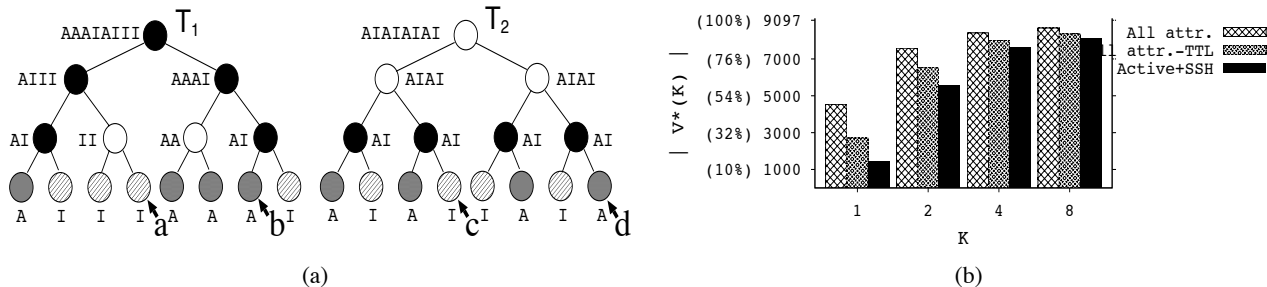


Figure 5: (a) Example of two fingerprint trees,  $T_1$  and  $T_2$ , for two distinct prefix-preserved anonymized traces; (b)  $K$ -vulnerability of trace *Trace-0618* in respect to Section 3 fingerprint attributes.

### 4.1 Calculating worst-case de-anonymization

To formalize the worst case de-anonymization, we wish to find the largest set of all anonymized addresses that are  $K$ -vulnerable when attacked using trace fingerprints  $F_t$ . We refer to this set as  $V^*(K)$ , and note that any set of  $K$ -vulnerable addresses obtainable by any de-anonymization technique using trace fingerprints  $F_t$ , (including the one presented by Coull et al. [22]) is a subset of  $V^*(K)$ . We present a very efficient algorithm that computes  $V^*(K)$  given  $F_t$ .

Recall that  $V(K, F_e, F_t, c)$  consists of the match sets for each anonymized host with size less than  $K$  and which contain the correct de-anonymization. For a fixed trace fingerprint  $F_t$ , let  $C(K)$  be the collection of all  $V(K, F_e, F_t, c)$  obtained by varying  $F_e$  and  $c$ . One can interpret  $C(K)$  as the collection of all attacking results using any combination of external fingerprints and cost functions.

Let elements in  $C(K)$  be partially ordered by the containment relation, i.e.,  $V \leq V'$  implies  $V \subseteq V'$ . We can show that for any value of  $K$  there exists an upper bound in  $C(K)$ , i.e.,  $\exists V^* \in C(K)$  such that  $V \subseteq V^*, \forall V \in C(K)$ . The following theorem states that  $V^*(K)$  exists and is related to the set of perfect external fingerprints  $F_e^* = F_t$ . For our theorem we need the following cost function

$$c^*(F_e(x), F_t(y)) = \begin{cases} 0 & \text{if } F_e(x) = F_t(y), \\ \epsilon & \text{for } \epsilon > 0, \text{ otherwise.} \end{cases} \tag{5}$$

The theorem states that  $C(K)$  is upper bounded by  $V^*(K) = V(K, F_e^*, F_t, c^*)$ . The proof of the following theorem can be found in Appendix D.

**Theorem 4.1.** *For an arbitrary set of trace and external fingerprints  $F_t, F_e$ , any cost function  $c$ , and any value of  $K$ :  $V(K, F_e, F_t, c) \subseteq V(K, F_e^*, F_t, c^*)$ , i.e.,  $V^*(K) = V(K, F_e^*, F_t, c^*)$  is an upper bound in  $C(K)$ .*

The calculation<sup>3</sup> of  $V^*(K)$ , and thus the evaluation of worst-case de-anonymization, is even more efficient than the execution of the attack described in Section 2. Because fingerprints are assumed perfect, cost based tree-matching is not needed. In fact, the entire computation can be performed on the trace fingerprint tree so that it is possible to compute  $V^*(K)$  in linear time ( $O(|\mathcal{A}|)$ ).

More specifically, Appendix E shows that the worst-case match sets for a host are determined by the number of “white” nodes in the fingerprint tree: the size of the match set for host  $y$  is  $2^{W(y)}$  where  $W(y)$  is the number of white nodes on the path from leaf  $y$  to the root of  $T(F_t)$ . As an example, Figure 5(a) shows trace fingerprint tree  $T_1$  with two marked leaves  $a$  and  $b$  and trace fingerprint tree  $T_2$  with two marked leaves  $c$  and  $d$  ( $a, b, c$ , and  $d$  are the leaves pointed by an arrow). Leaf  $a$  has one “white” vertex on its path to the root and thus  $|\beta(a, F_t, F_e^*, c^*)| = 2^1$ . Leaf  $b$  has no “white” vertex on its path to the root and thus  $|\beta(b, F_t, F_e^*, c^*)| = 2^0$ . Leaf  $c$  has two “white” vertices on its path to the root and thus  $|\beta(c, F_t, F_e^*, c^*)| = 2^2$ .

### 4.2 Worst-case $K$ -vulnerability under full prefix preservation

We can now apply the worst-case analysis to the trace studied experimentally in Section 3. Figure 5(b) shows the worst-case  $K$ -vulnerable hosts,  $|V^*(K)|$ , for  $K = 1, 2, 4, 8$  applied to *Trace-0618*. The bars labeled “All attr.”

<sup>3</sup>A Java source code of an algorithm that computes  $V^*(K)$  is available for download at <http://www-net.cs.umass.edu/~ribeiro/deanonymization/>

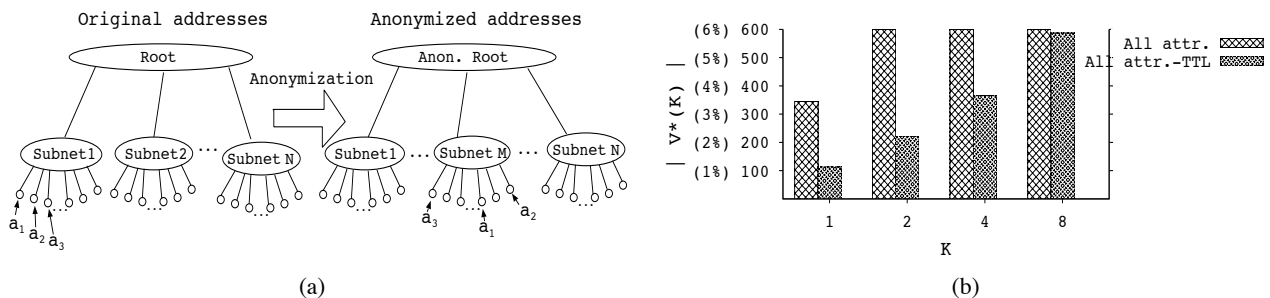


Figure 6: **(a)** Example of partial prefix preservation. The network administrator defines  $N$  subnets. Addresses sharing the same subnet in the original address space also shares the same subnet after the anonymization; **(b)**  $K$ -vulnerability of trace *Trace-0618* using partial prefix preservation in respect to Section 3 fingerprint attributes.

represent all attributes seen in Section 3. This graph shows that a well informed adversary can do significantly more damage to the anonymization function: 44% of the addresses are 1-vulnerable (uniquely identifiable), compared with 17% using external fingerprints obtained by probing the network.

Figure 5(b) also shows the worst case for other sets of attributes: “All attr. - TTL”, which includes all attributes except TTL, and “Active + SSH”, which includes only host activity and SSH traffic. It is clear from the graph that the TTL attribute only significantly increases the number of 1-vulnerable addresses. The “Active + SSH” bars show that other TCP port attributes besides “SSH” can significantly increase the precision of the de-anonymization. Figure 3(b) compares the “All attr.” upper bound with the results obtained with external fingerprints from *External-0618*. We can see that the perfect fingerprints perform at least two times better than the probed external fingerprints.

## 5 Analysis of partial prefix preservation

Given the vulnerability of traces to fingerprinting attacks like ours and others noted in the literature, it is natural to consider sacrificing the utility of the trace to increase privacy. We refer to such techniques as *partial prefix preservation*, because some prefix relationships are not preserved, replaced instead by unconstrained randomization of addresses or parts of addresses.

Pang et al. proposed a version of partial prefix preservation for internal addresses in their enterprise trace anonymization scheme [19]. There the subnet and host portions of an address are treated independently. Each subnet is comprised of a group of addresses that share the same  $32 - b$  most significant bits.<sup>4</sup> If two addresses share the first  $32 - b$  significant bits in the original trace then they also share the same first  $32 - b$  significant bits in the anonymized trace. Within subnets suffixes are not preserved; instead addresses are randomized within the subnet. Figure 6(a) illustrates this method. In this figure we label three hosts  $a_1$ ,  $a_2$ , and  $a_3$ . Note that  $a_1$ ,  $a_2$ , and  $a_3$  are just host labels, not the anonymized or unanonymized addresses of the hosts. Figure 6(a) shows both the original and the anonymized partial address-space trees. Note that  $a_1$ ,  $a_2$ , and  $a_3$  are consecutive addresses at the original address space. However, after the anonymization, the only address relation preserved is that  $a_1$ ,  $a_2$ , and  $a_3$  still belong to the same subnet. The subnet changed its address from “Subnet 1” to “Subnet M”.

**Worst-case analysis of partial prefix preservation** While this technique seems sensible, Pang et al. do not provide an analysis justifying this choice of partial prefix preservation. Attacks against this technique have since been demonstrated [22]. We now apply our analysis tools to evaluate the impact of partial prefix preservation on our sample trace.

The calculation of  $V^*$  under partial prefix preservation requires a modest variation of the technique described in the last section (details can be found in Appendix F). The results of the analysis appear in Figure 6(b) which plots  $|V^*(K)|$  against  $K = 1, 2, 4, 8$  for a subnet size of  $b = 8$  bits. Once the change in the y-axis has been noted, Figure 6(b) and Figure 5(b) (full prefix preservation) can be compared. We see that partial prefix preservation is much safer than full prefix preservation, reducing the unique de-anonymizations from about 4000 (44%) to 345 hosts (just under 4%). Partial prefix preservation also achieves much smaller values of  $|V^*(K)|$  for  $K = 2, 4, 8$ .

<sup>4</sup>In [19] the value of  $b$  is not necessarily constant among subnets; for the sake of simplicity we consider a constant value of  $b$  for all subnets.

Another interesting result comes from removing some attributes from the fingerprints. Figure 6(b) also plots the value of  $|V^*(K)|$  for attributes “All attr.-TTL”. Note that unlike full prefix preservation, the TTL attribute also significantly increases the number of 2 and 4-vulnerable addresses. Fingerprint attributes “Active+SSH” are not shown in Figure 6(b) as  $|V^*(1)| = |V^*(2)| = 0$  and  $|V^*(4)| = |V^*(8)| = 7$ . It is interesting to note that a fingerprint (“Active+SSH”) which is highly efficient for full prefix preservation is highly inefficient for partial prefix preservation. There is a simple explanation for the bad performance of fingerprint “Active+SSH” in partial prefix preservation: in our trace we found that most addresses with SSH traffic were clustered within a small number of subnets.

We also found that the size of the subnet, defined by  $b$ , can be varied with somewhat predictable results. When  $b = 4$  bits,  $|V^*(1)| = 1039$ , or 11% of hosts. When  $b$  is increased to 11 bits, few prefixes are preserved and  $|V^*(1)| = 106$ , or just 1%. Larger  $b$  improves host anonymity at the cost of host prefix information.

## Trading-off privacy and utility

Overall, for our trace, we find that partial prefix preservation as described by Pang et al. has a significant positive impact on trace anonymity. At the same time, 345 uniquely disclosed hosts may still be a concern to some publishers. In addition, these unique de-anonymizations confirm and help to explain the attack result found by Coull et al [22] in which a small number of distinctive servers are re-identified despite partial prefix preservation.

Interestingly, we found that it was possible to identify randomized subnets with very high likelihood. The concept of  $K$ -vulnerability can be applied not only to anonymized addresses but also to anonymized subnets. Consider the example in Figure 6(a). An adversary may be able to uniquely map anonymized “Subnet M” into its unanonymized counterpart “Subnet 1”. Note that this does not mean that an addresses inside anonymized “Subnet M” is also uniquely matched to an address inside unanonymized “Subnet 1”. In the following experiment we use  $b = 8$  bits and fingerprints “All attr.-TTL” obtained from *Trace-0618*. We choose fingerprint attributes “All attr.-TTL” instead of “All attr.” following Pang et al. [19] which removes TTL values from the trace. This experiment shows that 96% of all subnets, that have more than one active host, are 1-vulnerable and the remaining 4% are 2-vulnerable. When  $b = 4$  bits we find 17% of the subnets to be 2-vulnerable (of which 13% are 1-vulnerable). And when  $b = 11$  bits all subnets are 1-vulnerable. These results suggest that as subnets get larger their corresponding fingerprint are more likely to be unique.

In the light of these results we propose that for large enough values of  $b$  (such that most subnets 1-identifiable), one can increase trace utility (keeping the same upper bound anonymity level) by applying full prefix preservation from the subnet level up to the root.

Note finally that our upper bound also allows data publishers to reduce trace utility (e.g. removal of TCP port information for a given address) when it is absolutely necessary to ensure anonymity against a set of fingerprint attributes. This can be seen as an improvement over a previous recommendation [22] which advocates for the removal of TCP port numbers from all records in the trace. Our results also show that an attacker is able to achieve significantly more damage when address fingerprints contain TTL attributes. This result corroborates the notion in [19] that TTL attributes should not be published in the trace.

## 6 Conclusion and Future work

We have analyzed a novel attack on prefix-preserving trace anonymization that encompasses other proposed attacks, and can be adapted to partial prefix preservation and to alternative information sources. We demonstrated the effectiveness of this attack on a real trace, and measured experimentally the impact of the accuracy of information on the adversary’s success. Past works describing attacks on trace anonymization have left trace publishers without methods to evaluate the risks of publication for their traces and with few mitigation techniques. Our analysis techniques allow trace publishers to compute an upper bound for the risk of host de-anonymization in the context of adversaries assumed capable of collecting a given class of external information. In the future we hope to use these techniques to formally evaluate partial prefix preservation alternatives which can maximize utility relative to a desired level of trace privacy. We would also like to consider the application of this kind of attack to the external hosts in the trace. The challenges are that for external addresses structural information is less informative, and that probing must be limited to some known set of popular destinations since exhaustive probing of internet addresses would be infeasible.

## References

- [1] Y. Liu, X. Liu, L. Xiao, L. M. Ni, and X. Zhang, "Location-aware topology matching in P2P systems," *Proceeding of the IEEE INFOCOM 2004*, vol. 4, pp. 2220–2230, 2004.
- [2] M. Rajab, F. Monrose, and A. Terzis, "Fast and Evasive Attacks: Highlighting the challenges ahead," in *Proceedings of the 9th International Symposium on Recent Advances in Intrusion Detection (RAID)*, Hamburg, Germany, Sept. 2006.
- [3] "Tcprpriv," <http://ita.ee.lbl.gov/html/contrib/tcprpriv.html>.
- [4] R. Pang and V. Paxson, "A High-Level Programming Environment for Packet Trace Anonymization and Transformation," in *Proceedings of the ACM SIGCOMM Conference*, August 2003.
- [5] J. Xu, J. Fan, M. Ammar, and S. Moon, "Prefix-preserving IP address anonymization: Measurement-based security evaluation and a new cryptography-based scheme," in *Proceedings of the 10th IEEE International Conference on Network Protocols (ICNP'02)*, Paris, France, November 2002.
- [6] R. Ramaswamy and T. Wolf, "High-speed prefix-preserving IP address anonymization for passive measurement systems," *IEEE/ACM Transactions on Networking*, vol. 15, no. 1, January 2007.
- [7] "Gateway link measurements at umass amherst," <http://www-net.cs.umass.edu/dag/>.
- [8] "Enterprise tracing project," <http://www.icir.org/enterprise-tracing/>.
- [9] "The passive measurement and analysis project," <http://pma.nlanr.net/>.
- [10] "The skitter project," <http://www.caida.org/tools/measurement/skitter/>.
- [11] "The internet traffic archive," <http://ita.ee.lbl.gov/>, Apr. 2000.
- [12] "Network tools and traffic traces at university of napoli federico ii," <http://www.grid.unina.it/Traffic/>.
- [13] T. McGregor, H. Braun, and J. Brown, "The NLANR network analysis infrastructure," *IEEE Communications Magazine*, vol. 38, no. 5, pp. 122–128, May 2000.
- [14] K. Tai, "The tree-to-tree correction problem," *Journal of the Association for Computing Machinery (JACM)*, vol. 26, no. 3, pp. 422–433, 1979.
- [15] P. Bille, "A survey on tree edit distance and related problems," *Theor. Comput. Sci.*, vol. 337, no. 1-3, pp. 217–239, 2005.
- [16] A. Slagell and W. Yurcik, "Sharing Computer Network Logs for Security and Privacy: A Motivation for New Methodologies of Anonymization." in *Proceedings of SECOVAL: The Workshop on the Value of Security through Collaboration*, September 2005.
- [17] D. Koukis, S. Antonatos, and K. Anagnostakis, "On The Privacy Risks of Publishing Anonymized IP Network Traces," in *Proceedings of the 10th IFIP Open Conference on Communications and Multimedia Security*, October 2006.
- [18] T. Brekne, A. Årnes, and A. Øslebø, "Anonymization of IP Traffic Monitoring Data: Attacks on Two Prefix-preserving Anonymization Schemes and Some Proposed Remedies," in *Proceedings of the Workshop on Privacy Enhancing Technologies (PET 05)*, May 2005.
- [19] R. Pang, M. Allman, V. Paxson, and J. Lee, "The devil and packet trace anonymization," *SIGCOMM Comput. Commun. Rev.*, vol. 36, no. 1, pp. 29–38, 2006.
- [20] M. Z. et al., "<http://lcamtuf.coredump.cx/p0f.shtml>."

- [21] M. Neuhaus and H. Bunke, “Automatic learning of cost functions for graph edit distance.” *Information Sciences*, vol. 177, no. 1, pp. 239–247, 2007.
- [22] S. Coull, C. Wright, F. Monrose, M. Collins, and M. Reiter, “Playing devil’s advocate: Inferring sensitive information from anonymized network traces,” in *Proceedings of the Network and Distributed System Security Symposium*, 2007.
- [23] S. M. Bellovin, “A technique for counting NATed hosts,” in *IMW ’02: Proceedings of the 2nd ACM SIGCOMM Workshop on Internet measurement*, 2002, pp. 267–272.
- [24] B. Ribeiro, W. Chen, G. Miklau, and D. Towsley, “On the loss of privacy in packet trace anonymization,” Department of Computer Science, University of Massachusetts at Amherst, Tech. Rep. TR–07, Jun. 2007.

## APPENDIX

### A Obtaining external information address fingerprints

In what follows we compile a small list of possible external fingerprint attributes.

- *TCP port attribute (Web, FTP, ...)*. (Active probing) To test whether a address could have TCP activity on a given port one can simply try to use TCP to connect to that port.
- *Traffic volume*. Traffic volume can be divided into levels. For instance, all main network servers (main web servers, main e-mail servers and so on) may be inferred as the busiest ones for that type of service. The number of levels depends on the adversary’s knowledge of network servers traffic loads.
- *Flow sizes*. Many file transfers over the Internet have very specific flow size signatures.
- *Packet timing*. Internet activities (such as browsing the Web) can have easily identifiable packet inter-departure signatures.
- *Network changes*. Networks normally evolve with time. Such evolution can be available to the adversary. Abnormal traffic patterns or scheduled changes to the network that are known to the adversary creates the possibility of a new fingerprint attribute. Here even a harmless 4pm weekly department-wide coffee-break can result in a fingerprint attribute.

### B Sources of fingerprint attribute mismatches and their counter-measures

The following list is a short but representative compilation of the troubles that adversaries face:

- *Change in address fingerprints*. It is possible that some hosts/services are active (inactive) during the time of trace collection and inactive (active) when the adversary obtains address fingerprints. Under such circumstances it is likely that some address fingerprints appearing in the trace are distinct from their corresponding external information fingerprints. A good representative of this class are IP addresses allocated using the DHCP protocol (dynamic allocation).
- *Host or host service is active but had no traffic recorded*. Some hosts or host services may be active at the time of the trace collection but have no traffic recorded in the trace. This creates a fingerprint mismatch as the adversary sees a potential source of traffic that is not present in the trace.
- *Firewalls, IDS*. Firewalls and Intrusion Detection Systems (IDS) may limit the adversary’s ability to determine address fingerprints or even masquerade them.
- *NAT boxes*. Many distinct hosts can share an unique IP address which may induce contradictory fingerprints attributes at the same IP address.

*Attenuating fingerprint attribute mismatches.* Each item of the aforementioned list can be attenuated if the adversary takes the following precautions:

- *Dynamic IP address allocation.* An adversary that has access to a trace that spans over a long time frame may use it to obtain which portions of the address space may have temporal fingerprint changes. Information on which portions of the address space are more “dynamic” can actually help the adversary if the adversary has this type of external information (e.g., which portions of the network use DHCP). Another way to use such fingerprint is to mark the main e-mail server of a network as “stationary” and the rest of the e-mail servers as potentially “dynamic”.
- *Host or host service is active but had no traffic recorded.* A service or a host that sees has almost no recorded traffic in the trace is likely to .
- *Firewalls, IDS.* This is the hardest mismatch to cope with. Here the adversary cannot trust some fingerprint attributes and must resort to other types of external information, such as personal web-traffic preferences. For instance, Google.com flow sizes vary according to user language settings.
- *NAT boxes.* Might be detected in the trace [23]. If the adversary knows the true address of some NAT boxes, it can be used in the adversary’s advantage as a new “true, false or undefined” fingerprint attribute: “is it a NAT box?”.

## C Sample attack

Let  $T(F_t)$  and  $T(F_e)$  denote the trace and external fingerprint trees depicted in Figures 2(b) and 2(a), respectively. In this example we use the edit cost function  $c^*$  (equation (5)). Let’s follow the disclosure of address  $t_{15}$  of Figure 2(b). The algorithm works on the fingerprint trees in a deep-first-search-similar order, starting from the root of  $T(F_e)$ . It first decides which vertex of  $T(F_t)$ ,  $t_2$  or  $t_3$ , is the best match for vertex  $e_2$  of  $T(F_e)$ . Here, inner vertex fingerprints can help us decide which match is the most likely one. The fingerprint of  $e_2$  is **AAAI** and the fingerprints of  $t_2$  and  $t_3$  are **AAII** and **AAAI** respectively. The algorithm decides that  $t_3$  is the most likely match for  $e_2$ . The match  $(e_2, t_3)$  forces the match  $(e_3, t_2)$  due to the prefix preservation of the anonymization. The above illustrates our main optimization over the brute force search. Since  $f(e_2) = f(t_3)$ ,  $f(e_3) = f(t_2)$ ,  $f(e_2) \neq f(t_2)$ , and  $f(e_3) \neq f(t_3)$ , the minimum cost of the match  $(e_2, t_3)$  is zero rather than the minimum cost of  $(e_2, t_2)$  which is  $2\epsilon$  (for simplicity, we use  $(e, t)$  to denote that vertex  $e$  is de-anonymized or *matched* to vertex  $t$ ). If later in the process we believe for some reason that this decision was incorrect, we backtrack and follow the opposite match. But if the cost of matching the leaves of the  $(e_2, t_3)$  match is still zero (as predicted), then there is no need to verify the match  $(e_2, t_3)$ . The next step is to decide if the best match for  $e_4$  is  $t_6$  or  $t_7$ . Again we use fingerprints to decide the most likely match:  $(e_4, t_6)$ . Next we have to decide if  $e_8$  matches  $t_{12}$  or  $t_{13}$ . In this case we could go either way. The reason we cannot decide is that both  $t_{12}$  and  $t_{13}$  have the same labels. In this case we say that  $e_8$  is in the match set of  $t_{12}$  and  $t_{13}$ , or more formally:  $e_8 \in \beta(t_{12}, F_e, F_t, c^*)$  and  $e_8 \in \beta(t_{13}, F_e, F_t, c^*)$ . In the next step we find that  $t_{15}$  is a good match to  $e_{10}$  (thus  $e_{11}$  is a good match to  $t_{14}$ ). If we continue the algorithm we find that the edit distance is zero and that  $\beta(t_{15}, F_e, F_t, c^*) = \{e_{10}\}$ . After finishing the match  $(e_2, t_3)$ , we follow the match branch  $(e_3, t_2)$ . As  $f(e_4) = f(t_5)$  and  $f(e_4) \neq f(t_4)$ , we first follow the match  $(e_4, t_5)$ , from which we conclude that  $\beta(t_{10}, F_e, F_t, c^*) = \beta(t_{11}, F_e, F_t, c^*) = \{e_{14}, e_{15}\}$ . Likewise, we also find  $\beta(t_8, F_e, F_t, c^*) = \beta(t_9, F_e, F_t, c^*) = \{e_{12}, e_{13}\}$ . Note that we find all matches that have added costs (eq. 1) zero. The edit distance is therefore zero. However, when external fingerprints are imperfect, the edit distance may not be zero. More details of this algorithm can be found in [24].

## D Proof of the upper bound $V^*(K)$ (Theorem 4.1)

**Theorem** For an arbitrary set of trace and external fingerprints  $F_t, F_e$ , any cost function  $c$ , and any value of  $K$ :  $V(K, F_e, F_t, c) \subseteq V(K, F_e^*, F_t, c^*)$ , i.e.,  $V^*(K) = V(K, F_e^*, F_t, c^*)$  is an upper bound in  $C(K)$ .

*Proof.* In what follows we sacrifice formalism for clarity in order to convey the main idea behind the proof. We use the example in Figure 7 to illustrate the proof. In Figure 7 we represent four distinct address fingerprints.

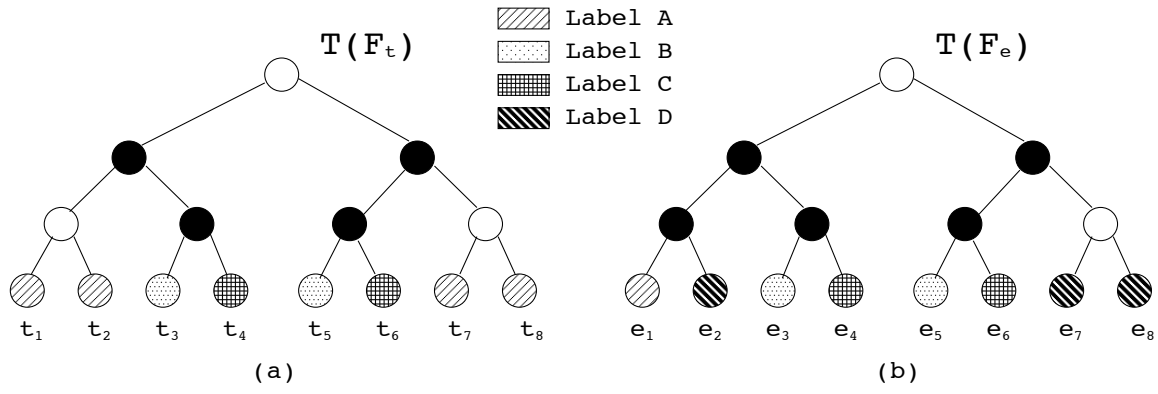


Figure 7: (a) Trace fingerprint tree; (b) External fingerprint tree.

Figure 7(a) and (b) show the trace and the external fingerprint trees, respectively, of our example. It is clear that  $t_a \in V^*(2)$ ,  $a = 3, 4, 5, 6$ ,  $t_i \notin V^*(2)$ ,  $i = 1, 2, 7, 8$ , and  $V^*(4)$  contains all addresses. Suppose that  $V^*(2) \subset V(2)$ . We show that the previous statement must be false. If the above is true then exists  $j \in \{1, 2, 7, 8\}$  such that  $t_j \in V(2)$ . Without loss of generality assume  $j = 8$ . Then  $|\beta(t_8, F_e, F_t, c)| \leq 2$ . But  $t_i \in V^*(2)$ ,  $i = 1, 2, 7, 8$  implies that any bijective mapping between  $t_8$  and  $t_1, t_2$ , or  $t_7$  is an isomorphism (i.e.,  $t_8$  is indistinguishable from  $t_1, t_2$ , or  $t_7$ ). Thus, for any value of  $m$ , the overall cost of matching  $(e_m, t_1)$  cannot be different than the overall cost of matching  $(e_m, t_8)$ . Thus  $e_m \in \beta(t_1, F_e, F_t, c)$  implies that  $e_m \in \beta(t_8, F_e, F_t, c)$ . The same is valid for  $t_2$  and  $t_7$ . From the above we conclude that  $\beta(t_8, F_e, F_t, c)$  must have at least four addresses and thus  $t_8 \notin V(2)$ .  $\square$

## E Upper bound computation, full prefix-preservation

**Lemma** Let  $T(F_t)$  be a trace fingerprint tree,  $y$  be a leaf of this tree, and  $N(y)$  be the set of all inner vertices from  $y$  (exclusive) to the root of  $T(F_t)$  (inclusive). Let  $W(y)$  be the number of white vertices in  $N(y)$ . Then  $|\beta(y, F_e^*, F_t, c^*)| \leq 2^{W(y)}$

*Proof.* If  $z \in N(y)$  is white, it means that the children of  $z$  are isomorphic. Thus there are  $2^{W(y)}$  leaves that can be mapped to  $y$  using an isomorphism.  $\square$

## F Upper bound computation, partial prefix-preservation

Let  $T(F_t)$  be a trace fingerprint tree,  $y$  be a leaf of this tree, and  $N(y)$  be the set of all vertices from  $y$  (exclusive) to the root of  $T(F_t)$  (inclusive). Let  $z$  be an inner vertex in  $T(F_t)$  and  $b(z)$  be the number of siblings of  $z$  that have the same fingerprint as  $z$ . Then  $W(y) = \sum_{z \in N(y)} b(z)$  and  $2^{W(y)}$  is number of leaves that can be mapped to  $y$  using an isomorphism.