

# The Limits of Traceability

Richard Clayton

University of Cambridge, Computer Laboratory, Gates Building, JJ Thompson Avenue,  
Cambridge CB3 0FD, United Kingdom

`richard.clayton@cl.cam.ac.uk`

**Abstract.** Traditional “traceability” (the flipside idea to “anonymity”) on the Internet attempts to identify the IP address that caused an action to occur. This is sufficient for an Internet Service Provider (ISP) to take action against the authorized user of that IP address. Law enforcement agencies usually need to go beyond this in order to identify the individual concerned. However, shared accounts, unavailable Caller Line Identification (CLI), and spoofing on Ethernets, mean that in the real world there is poor traceability for the “last hop”. This means that it is not possible to consider the information gathered as “conclusive evidence” suitable for a Court, but just as “intelligence”; an investigative tool, albeit a valuable one. Law enforcement officers should be especially wary that a sophisticated opponent might be able to frame an innocent bystander.

## 1. Introduction

It is often the case that an action has occurred on the Internet and one wishes to answer the question, “who did this?” The ability to track down the originator of an action is usually called “traceability” and it can be seen as the flipside idea to “anonymity”. It is of interest to regulators [1] and to law enforcement agencies (LEAs) [2] as well as to those who operate the Internet itself. There are a number of techniques available for establishing the originator of an action, many of which can be found described in the London Internet Exchange (LINX) “Best Practice” document [3].

The LINX document describes how to determine which IP address was the source of the action of interest and how to trace who was using that IP address at the relevant time. The document also considers what Best Practice procedures will ensure that this tracing can be performed as simply as possible; or, to put the same idea the other way around, which steps in the evidence chain are open to attack by a defense lawyer.

LEAs have seized upon traceability as a magic wand for the detection of “Hi-Tech” crime. Data retention (causing logs to be preserved for a known period) and data preservation (ensuring logs of special interest are not destroyed) have become hot topics at the inter-governmental level within the G8 and elsewhere. In the United Kingdom (UK) some LEA officials have gone so far as to propose data warehouses that will hold seven years worth of logging information from telephone companies and Internet Service Providers (ISPs) [4].

Traceability is usually talked about in the abstract without considering who is doing the tracing. However, the reason for doing the tracing cannot be ignored. Traceability begins to become far more problematic as one nears the actual source of the action and this poses very significant problems for a law enforcement official whereas an ISP may be able to ignore these complexities altogether.

Real world experiences, such as those of the author in the UK ISP industry, show that significant amounts of logging data are regularly collected and this can sometimes provide excellent levels of traceability. However, as soon as one wishes to identify individuals then one is reduced to deduction and inference rather than being able to be “sure”.

As an example, consider a dialup Internet access account that has been used for sending unsolicited bulk email (“spam”). When the evidence of this abuse arrives at the ISP the dialup account will be identified beyond question and it will then be disabled. No further investigation is needed. If, instead, the dialup account had been used for posting illegal material (for example, sexual images of children are unlawful in most jurisdictions) then law enforcement officers will be interested in exactly who was using the dialup account to post the material. This can be, as will be explained in detail below, extremely hard to determine.

In fact, it turns out that almost every method of accessing the Internet poses significant traceability problems over “the last hop”. This means that LEAs should cease to consider traceability as an evidential tool (proof that someone is guilty) and should only be considering it as an investigative tool (pointers to people who are worth paying more attention to). In particular, where the quarry is known to be technically sophisticated, careful attention must be paid to the possibility that an innocent person may be “framed” by the guilty and heavy-handed intervention may serve only to “tip off” the actual miscreant.

Typical problems with various access methods will now be described in turn.

## **2. Dialup Access (Modems and ISDN)**

Dialup access is one of the most common forms of access to the Internet and there are great many ISPs offering this type of connection. A communications link is made over ordinary voice telephone lines (analogue or digital) using a modem or ISDN terminal adapter. The session starts with the user authenticating themselves to the ISP, usually by a simple password, and they are then given an IP address and connected to the Internet.

There are two ways in which it is possible to tie a dialup connection to a particular individual. One can examine the credentials that were given when the dialup account was opened and one can use CLI (Caller Line Identification) to determine the source of any particular call within the telephone network.

Traditionally, one paid an ISP for a dialup account so the credentials given when accounts were opened would include some means of payment. The user would have to continue making payments at regular intervals or their account would lapse. It is most unusual for dialup accounts to be paid for in cash and the financial community can, upon production of suitable authority, trace the owner of a checkbook or credit card. If the credentials were forged then the ISP would discover this fairly promptly (usually when the credit card company refused payment) and the account will be closed – if indeed it was ever opened.

The recent invention of other business models based on telephone interconnect charges, advertising, or mere hot air, have meant that it is no longer necessary for the ISP to collect money from their subscribers directly. The ISP will still attempt to collect customer information (if only to profile their user base for the benefit of advertisers) but it has no way of establishing whether it is accurate.

With credentials becoming increasingly dubious as an identification mechanism, this leaves just the CLI information as a traceability tool. The CLI can be recorded by the ISP’s equipment and in principle it provides the phone number that made each individual call to the ISP. However, in the United Kingdom (UK), and doubtless elsewhere, it has significant limitations:

- the user may have requested the phone company to withhold CLI

Even if CLI is not suppressed for all calls, dialing a special code (141 in the UK) will disable CLI for an individual call. The CLI is still available to a telco in an “engineering” form, but not to an end system. UK ISPs are to be given access to “engineering CLI” once a suitable Code of Practice has been formulated, but this arrangement is not yet in place.

- CLI may be lost at telco boundaries

It is unusual for trans-Atlantic calls to have any CLI, and in the past CLI has failed to travel between various UK based networks with any reliability.

- the CLI provided may be generic

One of the side effects of using some discount phone schemes can be that the call appears to originate at the discount provider rather than at the true source. Therefore traceability will depend upon how much logging has been enabled within the third party system.

The CLI may belong to a company, with many users dialing out through a PABX. Traceability will depend upon what types of records are kept of usage by particular telephone extensions.

In order to prevent serial abuse UK ISPs who are providing “free” services usually require that CLI be presented before calls are accepted. When abuse has occurred no further accounts can be opened or operated by calls that originate from the same CLI. The ISPs who provide “paid” services tend not to require CLI since they would reckon to identify serial abusers from re-use, for payment, of the same credit card number.

Even where traceability appears to be provided, either by credentials or by CLI, LEAs have considerable further difficulties to overcome with both forms of identification. It may be straightforward to determine which account was being used but if the details of access to the account are widely known then the individual who used it may be impossible to trace. The password for a company account may be known to dozens of people any of whom could pass it on to a friend. Passwords are often inadvertently posted to support newsgroups when evidence of login problems are being discussed and of course the sysadmins at the ISP could well be in a position to collect the passwords of any user they wish.

### **3. New Technologies**

There are a number of relatively new technologies for providing access to the Internet.

High-speed access is becoming available over phone lines or cable networks. In general terms, the IP traffic is transferred across a lower level transport system (such as ATM). Standard traceability techniques, as described in the LINX document, will determine which subscriber’s account has been used at the ISP level. The equivalent of CLI in this type of scheme is that work will have been done at the transport level to set up the path between customer and ISP. Provided that suitable records have been kept it should be straightforward to determine the customer premises that were involved.

However, these records may not exist in a suitable form where the system automatically configures paths to ISPs using its own authentication scheme. If this is the case, as it is today in the UK for ADSL access via British Telecom’s infrastructure, then there still remains the possibility that account and password details have been leaked and an imposter must be located.

In passing, it should be noted that it is, in principle, possible to determine the origin of a dialup call even without the presence of CLI. The telco companies involved will have billing records that include every single call and with accurate time and duration information a particular call could be located. However, because of the high volume of calls made this would be a non-trivial task and it is not a service that the telcos currently make available. However, on lower volume systems, such as ADSL networks, it may be that appropriate types of record exist and can be searched cost-effectively.

Internet access from mobile (cell) phones is also available. This may involve plugging a laptop into the handset or, with the current WAP system and forthcoming 3GPP devices, the computer may be within the phone itself.

Dialup access to the Internet from a mobile handset has the same set of traceability issues as does dialup from fixed lines. Once again, when faced with abuse the ISP can block the account or the CLI whereas LEAs have the same set of problems in establishing the actual identity of the miscreant. There is, however, extra complexity for law enforcement when a “prepaid mobile” is used. In many countries there are no registration requirements for the ownership of “prepaid mobiles” where the call minutes are bought in advance and both phones and time can be purchased for cash. Therefore, even if the phone number is available, there will be no record of the identity of the owner. Technology does exist for determining the physical location of the phone, but the accuracy of this is limited. In practice, most traceability for prepaid mobiles is done by an analysis of billing records to try and establish the likely identity of the owner from the record of who they have been in contact with.

#### **4. Leased Lines**

A leased line is, in general, a permanently configured pathway across the telco network from customer to ISP. Although the term originally described a dedicated piece of copper, it would be rare for specific bits of cable to be involved except for the very last hop to the customer premises. The rest of the pathway will be virtual circuits running over shared hardware.

Because of the way in which the circuits are provided it might be assumed that the customer end is fixed and cannot be hijacked or impersonated, except by physical attacks on the cabling or by corruption within the telco itself. These possibilities cannot be entirely discounted as experience with burglar alarm circuits has demonstrated [5].

However, in the case of leased lines, a much more significant traceability issue arises for LEAs. They will need to determine which workstation at the customer end of a leased line is involved in the events that they wish to investigate. Of course this is potentially a problem with any of the Internet access methods that are described above as well. It would be most unusual to connect a large company network to the Internet over a mobile phone, but many companies will use ISDN or ADSL to access the Internet and the same type of issues will arise there as well.

The basic difficulty is that the logging records at the ISP will ensure that events can be traced to the leased line, but the records will not usually be able to distinguish between different machines at the customer’s premises.

There are two general schemes for allocating IP addresses to customer networks. The first is to allocate a subnet, a group of 2, 6, 14, 30, 62... addresses of which one is used for the customer “router” and the others are available for individual machines. The second scheme is called NAT (Network Address Translation) [6] where a single IP address will be visible to the Internet as a whole and “private” address space will be used on the customer network. A specialist machine on the customer premises will keep track of all open connections and will distribute incoming data to the appropriate customer workstation. It is usual, but not ubiquitous, for NAT to be used for dialup and subnets for leased lines.

In principle, traceability can be extended from the open Internet onto the private IP network so as to determine the workstation that is of interest. However, there are a number of practical difficulties with this and the first one is that the customer may not be keeping sufficient logging. If NAT is being used then a record would be needed of the connections it has handled – which would be a substantial amount of data and is very unlikely to have been recorded. In both the NAT case and the subnet case, there may not be any records of which machines were using particular IP addresses. It is not unusual for addresses to be dynamically allocated by the DHCP protocol [7] and the server may not keep particularly useful records. In all cases where records are being kept there should be some caution in the minds of law enforcement officers as to whether these records are in any sense trustworthy or indeed whether the mere attempt to inspect them will “tip off” the miscreant about the investigation.

It is not only the lack of historical records that makes traceability problematic on customer networks. Even when records exist and even when one is tracking someone down in real-time rather than trying to untangle a previously occurring event there is a further problem. A technically competent person (or one with access to sophisticated tools) may be able to “hide” themselves on an Ethernet network in such a way as to incriminate an innocent user. This poses a problem for law enforcement officers in that a high profile “raid” on the innocent may serve to “tip off” someone nearby.

In order to understand how to “hide” on a Ethernet (which is by far the most common local area network technology in use today) one must remember that it is fundamentally a broadcast medium. Traffic will arrive from the Internet down the single wire from the ISP but is then broadcast to all the machines on the Ethernet. In the Ethernet environment the addressing is done by “MAC address” and network interfaces will only handle traffic for their particular address. Protocols such as ARP are used to associate IP address and MAC address together.

MAC addresses are intended to be unique and fixed, serving as a machine identifier. The IP addresses may be statically allocated by an administrator, in which case traceability is achieved by consulting the administrator’s records. Alternatively, it may be dynamically allocated by a protocol such as DHCP, in which case inspection of the DHCP server records will provide a mapping to a MAC address and possibly also an indication of information such as the Windows NT domain identifier being used.

Where historic records are unavailable (or cannot be trusted to be accurate) traceability might be achieved by real-time monitoring of the network if the target is still using it. However, even if the target is still using the network, they may be hard to physically locate. Time Domain Reflectometers are a classic way of locating where devices are attached, and can be used to locate extra kit that should not be connected to the network. However, in reality their output can be difficult to interpret, and switches and hubs will segment the network so that the search will take a long time. They will of course be of no use at all where you know that you are seeking an authorized attachment that is behaving in an unauthorized manner.

It is good practice to keep a record of the MAC addresses that belong to particular pieces of hardware to assist in tracing network packets to their source. However, even though the addresses are usually still stored in the hardware, modern network drivers tend to handle the MAC addresses within the software. This means that they are trivially changed to any value required – possibly “borrowing” other people’s settings. Difficulties arise if the machine the MAC address is borrowed from is still running on the network, but if it is switched off (or subject to a disabling denial of service attack) then it can be readily impersonated.

It is also possible to use “spare” IP addresses or to hijack another machines’ IP address. Again, complications arise if the IP address is simultaneously in use by its real owner, because the ARP protocol will be aware that unusual events are occurring. This can be overcome by careful forwarding of data combined with gratuitous ARP traffic to undo anything that the real owner of the IP address may do. With suitable tools such as monitoring the ARP table or network “sniffers” [8] it is possible to determine that unusual things are occurring, but even if this is known, it will still be hard to identify the node that is misbehaving.

## **5. Using Anonymous Machines**

Since tracing can sometimes be effective, some people use a machine other than their own, so that it will not matter that it is traced.

“Cybercafes” are often cited as the ultimate in anonymous access to the Internet. They have all the same characteristics as a corporate network except that there is no formal record kept of the identity of those who are using particular machines. There will of course be informal records in that by putting in a personal appearance a risk is run of being remembered by the counter staff or recorded on CCTV in the street outside.

Cybercafes do have an interest in preventing the type of abuse that would cause their connectivity provider to close them down (such as allowing the bulk sending of unsolicited email). As such, there must be an expectation that some types of traffic (such as connections to the SMTP port 25) may be monitored in real time. As such it may well not be a “safer” environment for the miscreant than a university or company network.

Cybercafes are not unique in providing public access points. One can often find machines in libraries or hotel lobbies, where it is likely that the technical expertise of the system administrators will be somewhat less than in an organization whose main business is the provision of connectivity.

Some people (often called “hackers”) may use a machine without permission, remotely operating it with the intent of providing a “cut-off” so that tracing will have to be repeated from that point to locate the actual user. If the events being investigated took place in the past and the local records have been competently removed (or altered) then the trail is effectively dead. If traffic is ongoing then, with co-operation from the owner of the machine or their network administrator, it should be a relatively simple task to determine the IP address of the controlling machine – and this can be located (or not) using standard techniques. However, real cases such as the ROME Labs intrusion [9] show that it may be very difficult to gather this type of evidence in a form that can be confidently expected to survive a challenge in court.

## 6. Conclusions

Traditional traceability techniques allow one to determine the administrative owner of an IP address. This is of great practical value to ISPs and other network operators because they can then use the threat of disconnection from the Internet as a way of ensuring co-operation and future good behavior. Connectivity providers have almost always seen removal of connectivity as sufficient sanction in itself and have seldom seen value in identifying people so that they might be called to account in the “Real World”. It is therefore hardly surprising to find that the traceability that has been developed over the years works extremely well in determining where to “pull the plug” but is poor at going beyond that.

LEAs have a need to identify people. The current emphasis by policy makers on the retention and preservation of logging data by telcos and ISPs will not, in practice, deliver what law enforcement needs to do this identification. The likely result of current initiatives is to demonstrate that in many common circumstances traceability information is unavailable or incomplete. At present, there is little evidence that improving the quality of the logging will be tackled in parallel to increasing the quantity.

Practical traceability often comes from imaginative use of the available records. This will usually provide LEAs with information that will materially assist in determining who should be investigated. However, the circumstantial nature of the evidence that traceability depends upon means that it will often be unsuitable for use before a court. Where the information is ambiguous, it may lead the investigator along false trails that are close enough to the real target to “tip them off”. Where individuals are actively trying to hide, they may succeed in making innocent bystanders appear responsible for actions they never performed.

Although logging is a “hot topic” for privacy campaigners, this is partly because it is ascribed magical powers to track everyone, everywhere. The failure of traceability at the edges of the network has seldom been tackled or described. ISPs have ignored the issue because their needs for traceability stop at the account level and these needs are satisfied by the status quo. LEAs are still learning what is possible in cyberspace and are currently concentrating on their perception that what is important is the failure of traceability in the center of the network because of unsuitable log retention policies. The academic security community shies away from topics that “everyone knows”, even though everyone doesn’t, and has mainly concentrated on protecting systems from harm, rather than catching the bad guys after the event. Meanwhile it is “hackers”, using individual techniques in an ad hoc manner, that have provided the majority of the documentation of viable methods of avoiding being traced.

It is now time to pull these disparate approaches together and understand more clearly what traceability can actually deliver and where its limits may lie. This paper has been a first attempt at performing this important task.

## References

1. Coffee, J.C.Jr.: Brave New World? The Impact(s) of the Internet on Modern Securities Regulation. *The Business Lawyer* 52(4) (1997).
2. President's Working Group on Unlawful Conduct on the Internet: The Electronic Frontier: The Challenge of Unlawful Conduct Involving the Use of the Internet. US Department of Justice (March 2000)
3. LINX Content Regulation Committee, Clayton et al: LINX Best Current Practice – Traceability. Version 1.0. (18 May 1999) <http://www.linx.net/noncore/bcp/traceability-bcp.html>
4. Gaspar, R.: NCIS submission on communications data retention law. <http://www.cryptome.org/ncis-carnivore.htm>
5. Anderson, R.: *Security Engineering: A Guide to Building Dependable Distributed Systems*. Wiley 2001. pp 213-216.
6. Egevang, K.: The IP Network Address Translator (NAT). Request for Comments 1631 (May 1994) <http://www.ietf.org/rfc/rfc1631.txt>
7. Droms, R.: Dynamic Host Configuration Protocol. Request for Comments 2131 (March 1997) <http://www.ietf.org/rfc/rfc2131.txt>
8. Graham, R.: Sniffing (Network wiretap, sniffer) FAQ. v0.3.3 (September 14, 2000) <http://www.robertgraham.com/pubs/sniffing-faq.html>
9. Sommer, P.: Intrusion Detection Systems as Evidence. RAID 98 (1998) [http://www.zurich.ibm.com/pub/Other/RAID/Prog\\_RAID98/Full\\_Papers/Sommer\\_text.pdf](http://www.zurich.ibm.com/pub/Other/RAID/Prog_RAID98/Full_Papers/Sommer_text.pdf)