

On The Limits of Steganography

Ross J. Anderson, Fabien A.P. Petitcolas

Abstract— In this paper, we clarify what steganography is and what it can do. We contrast it with the related disciplines of cryptography and traffic security, present a unified terminology agreed at the first international workshop on the subject, and outline a number of approaches—many of them developed to hide encrypted copyright marks or serial numbers in digital audio or video. We then present a number of attacks, some new, on such information hiding schemes. This leads to a discussion of the formidable obstacles that lie in the way of a general theory of information hiding systems (in the sense that Shannon gave us a general theory of secrecy systems). However, theoretical considerations lead to ideas of practical value, such as the use of parity checks to amplify covertness and provide public key steganography. Finally, we show that public key information hiding systems exist, and are not necessarily constrained to the case where the warden is passive.

Keywords— Cryptography, Copyright protection, Data compression, Image registration, Jitter, Motion pictures, Multimedia systems, Music, Observability, Pseudonoise coded communication, Redundancy, Spread spectrum communication, Software protection

I. INTRODUCTION

While classical cryptography is about concealing the content of messages, steganography is about concealing their existence. It goes back to antiquity: Herodotus relates how the Greeks received warning of Xerxes' hostile intentions from a message underneath the wax of a writing tablet, and describes a trick of dotting successive letters in a covertext with secret ink, due to Aeneas the Tactician. Kahn tells of a classical Chinese practice of embedding a code ideogram at a prearranged place in a dispatch; the same idea arose in medieval Europe with grille systems, in which a paper or wooden template would be placed over a seemingly innocuous text, highlighting an embedded secret message.

Such systems only make sense where there is an opponent. This opponent may be passive, and merely observe the traffic, or he may be active and modify it. A famous case dates back to 1586, when Mary Queen of Scots was conspiring to have Queen Elizabeth of England assassinated, with a view to taking over the English throne. However the cipher she used was broken, and the English secret police obtained the would-be assassins' names by forging a postscript to a letter she wrote to the chief conspirator, asking for "the names and qualities of the six gentlemen which are to accomplish the designation." This led to their arrest and execution, as indeed to Mary's the following year. In this century, postal censors have deleted lovers' X's from letters, shifted the hands of watches in shipments, and even rephrased telegrams; in one case, a censor changed "father

Manuscript received March 1997; revised August 1997. The work of F. Petitcolas was supported by the Intel Corporation under the Grant "Robustness of Information Hiding Systems."

The authors are with the University of Cambridge Computer Laboratory, Cambridge CB2 3QG, UK. E-mail (e-mail: rja14@c1.cam.ac.uk; fapp2@c1.cam.ac.uk).

is dead" to "father is deceased," which elicited the reply "is father dead or deceased?" [24].

The study of this subject in the scientific literature may be traced to Simmons, who in 1983 formulated it as the "Prisoners' Problem" [44]. In this scenario, Alice and Bob are in jail, and wish to hatch an escape plan; all their communications pass through the warden, Willie; and if Willie detects any encrypted messages, he will frustrate their plan by throwing them into solitary confinement. So they must find some way of hiding their ciphertext in an innocuous looking covertext. As in the related field of cryptography, we assume that the mechanism in use is known to the warden, and so the security must depend solely on a secret key that Alice and Bob have somehow managed to share.

There are many real life applications of steganography. Apparently, during the 1980's, Margaret Thatcher became so irritated at press leaks of cabinet documents that she had the word processors programmed to encode their identity in the word spacing, so that disloyal ministers could be traced. Similar techniques are now undergoing trials in an electronic publishing project, with a view to hiding copyright messages and serial numbers in documents [31].

Simmons' formulation of the Prisoners' Problem was itself an instance of information hiding. It was a ruse to get the academic community to pay attention to a number of issues that had arisen in a critical but at that time classified application—the verification of nuclear arms control treaties. The US and the USSR wanted to place sensors in each others' nuclear facilities that would transmit certain information (such as the number of missiles) but not reveal other kinds of information (such as their location). This forced a careful study of the ways in which one country's equipment might smuggle forbidden data past the other country's monitoring facilities [45], [47].

Steganography must not be confused with cryptography, where we transform the message so as to make its meaning obscure to a person who intercepts it. Such protection is often not enough. The detection of enciphered message traffic between a soldier and a hostile government, or between a known drug-smuggler and someone not yet under suspicion, has obvious implications; and recently, a UK police force concerned about criminal monitoring of police radios has discovered that it is not enough to simply encipher the traffic, as criminals detect, and react to, the presence of encrypted communications nearby [50].

In some applications, it is enough to hide the identity of either the sender or the recipient of the message, rather than its very existence. Criminals often find it sufficient for the initiator of a telephone call to be anonymous. Indeed, the main practical problem facing law enforcement and intelligence agencies is "traffic selection"—deciding which calls to intercept—and because of the huge volume of traf-

fic, this must usually be done in real time [29].

The techniques criminals use to thwart law enforcement vary from country to country. US villains use “tumblers”—cellular phones that continually change their identity, using genuine identities that have either been guessed or intercepted; in France, drug dealers drive around with a cordless phone handset until a dial tone is found, then stop to make a call [27]; while in one UK case, a drug dealer physically tapped into a neighbour’s phone [14]. All these techniques also involve theft of service, whether from the phone company or from one of its customers; so this is one field where the customer’s interest in strong authentication and the police interest in signals intelligence coincide. However authentication itself is not a panacea. The introduction of GSM, with its strong authentication mechanisms, has led crooks to buy GSM mobile phones using stolen credit cards, use them for a few weeks, and then dispose of them [55].

Military organisations also use unobtrusive communications. Their preferred mechanisms include spread spectrum and meteor scatter radio [40], which can give various combinations of resistance to detection, direction finding and jamming; they are vital for battlefield communications, where radio operators who are located are at risk of being attacked. On the Internet, anonymous remailers can be used to hide the origin of an email message, and analogous services are being developed for other protocols such as `ftp` and `http` [8], [17], [39].

Techniques for concealing meta-information about a message, such as its existence, duration, sender and receivers are collectively known as traffic security. Steganography is often considered to be a proper subset of this discipline rather than being co-extensive with it, so we shall now try to tie down a definition.

II. WHAT IS STEGANOGRAPHY?

Classical steganography concerns itself with ways of embedding a secret message (which might be a copyright mark, or a covert communication, or a serial number) in a cover message (such as a video film, an audio recording, or computer code). The embedding is typically parametrised by a key; without knowledge of this key (or a related one) it is difficult for a third party to detect or remove the embedded material. Once the cover object has material embedded in it, it is called a stego object. Thus, for example, we might embed a mark in a covertext giving a stegotext; or embed a text in a cover image giving a stego-image; and so on. (This terminology was agreed at the First International Workshop on Information Hiding [36]).

There has been a rapid growth of interest in this subject over the last two years, and for two main reasons. Firstly, the publishing and broadcasting industries have become interested in techniques for hiding encrypted copyright marks and serial numbers in digital films, audio recordings, books and multimedia products; an appreciation of new market opportunities created by digital distribution is coupled with a fear that digital works could be too easy to copy. Secondly, moves by various governments to restrict the availability of encryption services have motivated people

to study methods by which private messages can be embedded in seemingly innocuous cover messages. The ease with which this can be done may be an argument against imposing restrictions [16].

Other applications for steganography include the automatic monitoring of radio advertisements, where it would be convenient to have an automated system to verify that adverts are played as contracted; indexing of videicemail, where we may want to embed comments in the content; and medical safety, where current image formats such as DICOM separate image data from the text (such as the patient’s name, date and physician), with the result that the link between image and patient occasionally gets mangled by protocol converters. Thus embedding the patient’s name in the image could be a useful safety measure.

Where the application involves the protection of intellectual property, we may distinguish between watermarking and fingerprinting. In the former, all the instances of an object are marked in the same way, and the object of the exercise is either to signal that an object should not be copied, or to prove ownership in a later dispute. One may think of a watermark as one or more copyright marks that are hidden in the content.

With fingerprinting, on the other hand, separate marks are embedded in the copies of the object that are supplied to different customers. The effect is somewhat like a hidden serial number: it enables the intellectual property owner to identify customers who break their license agreement by supplying the property to third parties. In one system we developed, a specially designed cipher enables an intellectual property owner to encrypt a film soundtrack or audio recording for broadcast, and issue each of his subscribers with a slightly different key; these slight variations cause imperceptible errors in the audio decrypted using that key, and the errors identify the customer. The system also has the property that more than four customers have to collude in order to completely remove all the evidence identifying them from either the keys in their possession or the audio that they decrypt [6].

Using such a system, a subscriber to a music channel who posted audio tracks to the Internet, or who published his personal decryption key there, could be rapidly identified. The content owner could then either prosecute him, revoke his key, or both.

There is a significant difference between classical steganography, as modelled in the Prisoners’ Problem, and copyright marking. In the former, a successful attack consists of the warden’s observing that a given object is marked. In the second, all the participants in the scheme may be aware that marks are in use—so some effects of the marks may be observable (marks should remain below the perceptual threshold, but they may alter the content’s statistics in easily measurable ways). So a successful attack does not mean detecting a mark, but rendering it useless. This could be done by removing it, or by adding many more marks to prevent a court telling which one was genuine. Blocking such attacks may involve embedding a signature by the customer in the content [37] or involving a public

timestamping service in the marking process.

III. THE STATE OF THE ART

Prudent cryptographic practice assumes that the method used to encipher data is known to the opponent, and that security must lie in the choice of key. This principle was first enunciated by Kerckhoffs in 1883 [25], and has been borne out by long and hard experience since [24]. It should be an obvious requirement for protection mechanisms designed to provide evidence, as one can expect that them to be scrutinised by hostile expert witnesses in open court [2].

So one might expect that designers of copyright marking systems would publish the mechanisms they use, and rely on the secrecy of the keys employed. Sadly, this is not the case; many purveyors of such systems keep their mechanisms subject to non-disclosure agreements, sometimes offering the rationale that a patent is pending. So we will briefly survey a few systems that have been described in public, or of which we have information.

A. Simple systems

A number of computer programs are available that will embed information in an image. Some of them just set the least significant bits of the image pixels to the bits of the embedded information [53]. Information embedded in this way may be invisible to the human eye [28] but is trivial for an alert third party to detect and remove.

Slightly better systems assume that both sender and receiver share a secret key and use a conventional cryptographic keystream generator [41] to expand this into a long pseudo-random keystream. The keystream is then used to select pixels or sound samples in which the bits of the ciphertext are embedded [16].

Not every pixel may be suitable for encoding ciphertext: changes to pixels in large fields of monochrome colour, or that lie on sharply defined boundaries, might be visible. So some systems have an algorithm that determines whether a candidate pixel can be used by checking that the variance in luminosity of the surrounding pixels is neither very high (as on a boundary) nor very low (as in a monochrome field). Wherever a pixel passes this test, we can tweak its least significant bit to embed a bit of our message.

Such schemes can be destroyed in a number of ways by an opponent who can modify the stego-image. For example, almost any trivial filtering process will change the value of many of the least significant bits. One possible countermeasure is to use redundancy: either apply an error correcting code, or simply embed the mark a large number of times. For example, the “Patchwork” algorithm of Bender *et al.* hides a bit of data in an image by increasing the variance in luminosity of a large number of pseudorandomly chosen pixel pairs [10]; and a similar system was proposed by Pitas [38]. Much the same techniques can be used to mark digital audio as well.

One way in which we have attacked such systems is to break up the synchronisation needed to locate the samples in which the mark is hidden: pictures, for example, can be cropped. In the case of audio, we have developed a simple

but effective desynchronisation attack: we randomly delete a small proportion of sound samples, and duplicate a similar number of others. This introduces a jitter of a few tens of microseconds, which is tiny compared to the precision with which the original sounds were in most cases generated but is sufficient to confuse a typical marking scheme.

With a pure tone, we can delete or duplicate one sample in 8,000, and with classical music we can delete or duplicate one sample in 500, without the results being perceptible either to us or to laboratory colleagues. Using more sophisticated resampling and filtering algorithms, we can obtain a 1 in 500 jitter in pure tone, and 1 in 50 in speech, without making a perceptible difference. (The result for classical music can also be improved significantly, but the precise figure depends on the music.)

B. Operating in a transform space

A systematic problem with the kind of scheme described above is that those bits in which one can safely embed covert data are by definition redundant—in that the attacker will be unaware that they have been altered—and it follows that they might be removed by an efficient compression scheme. The interaction between compression and steganography is a recurring thread in the literature.

Where we know in advance what compression scheme will be used, we can often tailor an embedding method to get a quite reasonable result. For example, with .gif files one can swap colours for similar colours (those that are adjacent in the current palette) [23], while if we want to embed a message in a file that may be subjected to JPEG compression and filtering, we can embed it in multiple locations [26], [30] or, better still, embed it in the frequency domain by altering components of the image’s discrete cosine transform. A particularly detailed description of such a technique may be found in [13]; this technique, being additive, has the property that if several marks are introduced in succession, then they can all be detected (thus it is prudent for the originator of the content to use a digital timestamping service [51] in conjunction with the marking system, so that the priority of the genuine mark can be established). Other schemes of this kind include, for example, [11] and [26].

Such “spread spectrum” techniques are often tuned to the characteristics of the cover material. For example, one system marks audio in a way that exploits the masking properties of the human auditory system [12].

Masking is a phenomenon in which one sound interferes with our perception of another sound. Frequency masking occurs when two tones which are close in frequency are played at the same time. The louder tone will mask the quieter [21], [35]. However this does not occur when the tones are far apart in frequency. It has also been found that when a pure tone is masked by wideband noise, only a small band centred about the tone contributes to the masking effect [32]. Similarly, temporal masking occurs when a low-level signal is played immediately before or after a stronger one. For instance after we hear a loud sound, it takes a little while before we can hear a quiet one.

MPEG audio compression techniques exploit these characteristics [1], but it remains possible to exploit them further by inserting marks that are just above the truncation threshold of MPEG but still below the threshold of perception [12]. In general, a copyright mark's existence may be detected by statistical tests while it remains undetectable by humans; the real question is whether it can be damaged beyond later recognition without introducing perceptible distortion.

Embedding data in transformed content is not restricted to the 'obvious' transforms that are widely used for compression, such as discrete cosine, wavelet and fractal transforms. A recent interesting example has been suggested in [18]: this "echo hiding" technique marks audio signals by adding an echo. This echo might have a delay of 0.5 ms to signal a "0" and 1.0 ms to signal a "1;" these delays are too short to be perceptible in most circumstances but can be detected using cepstral transforms.

C. A general model

The general model of steganography we have developed in the above sections is that Alice embeds information by first applying a transform to the covertext, and then tweaking a subset of the bits of the transformed object that are now redundant. In this context, redundant means that a nontrivial subset of them, which is selected randomly to be of a given size, can have their values altered without this being detected easily or at all by an opponent who does not know which subset to examine.

We will not expect to find high bandwidth channels, as these would correspond to redundancy that could economically be removed. However, the design of compression schemes is limited in most cases by economic factors; the amount of computation that we are prepared to do in order to replace a certain amount of communication is not infinite. If we are prepared to do a little more work than the "normal" user of the system, we will be able to exploit a number of low-bandwidth stego channels.

However, the warden may be willing to do even more work, and the apparent redundancy which we exploit will fall within his ability to model. This may be especially so if the warden is a person with access to future technology—for example, a pirate seeking to remove the watermark or fingerprint embedded in a 1997 music recording using the technology available in 2047. This is a serious concern with copyright, which may subsist for a long time (typically 70 years after the author's death for text and 50 years for audio). Even where we are concerned only with the immediate future, the industry experience is that it is a "wrong idea that high technology serves as a barrier to piracy or copyright theft; one should never underestimate the technical capability of copyright thieves" [19]. Such experience is emphasised by the recent success of criminals in cloning the smartcards used to control access to satellite TV systems [5].

When such concerns arise in cryptography—for example, protecting traffic that might identify an agent living under deep cover in a foreign country—the standard solution is to

use a one-time pad; Shannon provided us with a proof that such systems are secure regardless of the computational power of the opponent [43]. Simmons provided us with a comparable theory of authentication, that has been applied in nuclear weapons command and control [46]. Yet we still have no comparable theory of steganography.

In the next section, we will discuss the formidable obstacles to such a theory, and indicate how some theoretical ideas have nonetheless led to useful improvements in the state of the art.

IV. THEORETICAL LIMITS

Can we get a scheme that gives unconditional covertness, in the sense that the one-time pad provides unconditional secrecy?

Suppose that Alice uses an uncompressed digital video signal as the covertext, and then encodes ciphertext at a very low rate. For example, the k th bit of ciphertext becomes the least significant bit of one of the pixels of the k th frame of video, with the choice of pixel being specified by the k th word of a shared one time pad. Then we intuitively expect that attacks will be impossible: the ciphertext will be completely swamped in the covertext's intrinsic noise. Is there any way this intuitive result could be proved?

We must first ask what a proof of covertness would look like. A working definition of a secure stegosystem might be one for which the warden cannot differentiate between raw covertext and the stegotext containing embedded information, unless he has knowledge of the key. As with cryptography, we might take the warden to be a probabilistic polynomial Turing machine in the case where we require computational security, and assume that he can examine all possible keys in the case where we require unconditional security.

In the latter case, he will see the actual embedded message, so the system must generate enough plausible embedded messages from any given stegotext, and the number of such messages must not vary in any usable way between the stegotext and a wholly innocent covertext.

This much is straightforward, but what makes the case of steganography more difficult than secrecy or authenticity is that we are critically dependent on our model of the covertext.

A. What if perfect compression existed?

Workers in information theory often assume that any information source can be compressed until there is no redundancy left. This assumption may be very useful in proving asymptotic bounds and capacity results, but has a rather curious effect when applied to steganography.

Suppose that such a perfect coding scheme were actually instantiated in a physical black box that could both compress and decompress data of a particular type (audio, video, whatever). Completely efficient compression means that the compressed objects would be dense in the set of bit strings of the same length. Thus Alice could take an arbitrary ciphertext message that she wants to hide and

run it through the decompressor. The result would be an acceptable audio recording, video film or whatever.

The above is not a rigorous proof. It is conceivable, for example, that a device might decompress a random bit string of length n to a particular type of object with a probability polynomial in $1/n$. This would suffice for many information theoretic results to go over, while invalidating the above argument. Nonetheless, it indicates that many classical intuitions of information theory serve us poorly when dealing with steganographic systems. It points to some interesting research problems in closing the gap between the two, and tells us that practical steganography is only an issue where compression is inefficient. Where efficient compression is available, information hiding will usually be either trivial or impossible, depending on the context.

B. Entropy

Entropy arguments are used in conventional information theory; how far will they get us in steganography?

Assuming that the material to be embedded is indistinguishable from random data (as would be the case were it competently encrypted), then entropy will be strictly additive: the entropy of the stegotext S will equal the entropy of the covertext C plus the entropy of the embedded material E :

$$H(S) = H(C) + H(E) \quad (1)$$

Thus in order to make our embedding process secure against an opponent who merely has to detect the presence or absence of embedded material, it appears that we have two alternatives:

1. Keep $H(E)$ much less than the uncertainty in the opponent's measurement of $H(C)$
2. Find some way of processing C to reduce its entropy by an amount that can then be made up by adding E . For example, one might use a noise reduction or lossy compression algorithm to remove some unnecessary information from C before embedding E .

The problem is that we do not know how competent our opponent is at measuring the entropy of the covertext we are using, or, equivalently, at discriminating signal from noise. We will often be up against an opponent of unpredictable power (a pirate attacking our system a generation from now); and these are precisely the circumstances where we may want a security proof.

But the more stegotext we give the warden, the better he may be able to estimate the statistics of the underlying covertext, and so the smaller the rate at which Alice will be able to tweak bits safely. The rate might even tend to zero, as was noted in the context of covert channels in operating systems [33]. However, as a matter of empirical fact, there do exist channels in which ciphertext can be inserted at a positive rate [16], and people have investigated correlations in various types of content such as digital video [48]. So measuring entropy may be useful in a number of applications.

But is there any prospect of developing steganographic techniques which we can prove will resist an opponent of arbitrary ability?

C. Selection channel

Our next idea is inspired by the correction channel that Shannon uses to prove his second coding theorem. (This is the channel which someone who can see both the transmitted and received signals uses to tell the receiver which bits to tweak; it produces various noise and error correction bounds [42].)

In a similar way, when Alice and Bob use a shared one-time pad to decide which covertext bit will be marked with the next ciphertext bit, we can think of the pad as a selection channel. If Willie is computationally unbounded, he can try all possible pads (including the right one), so the number of them which yield a plausible ciphertext must be large enough that he cannot reasonably accuse Alice of sending stegotext rather than an innocent message.

It may be useful at this point to recall the book cipher. The sender and receiver share a book and encipher a message as a series of pointers to words. So the cipher group "78216" might mean page 78, paragraph 2 and the 16th word. Book codes can be secure provided that the attacker does not know which book is in use, and care is taken not to reuse a word (or a word close enough to it) [24]. The book cipher is a kind of selection channel. The model of computation may appear to be different, in that with a book cipher we start off with the book and then generate the ciphertext, whereas in a stegosystem, we start off with the text to be embedded and then create the stegotext; but in the case where the selection channel is truly random (a one-time pad), they are the same, in that an arbitrary message can be embedded in an arbitrary covertext of sufficient length.

A repetitive book will have a lower capacity, as we will be able to use a smaller percentage of its words before correlation attacks from the context become possible. Similarly, if the covertext to be used in a stegosystem has unusual statistics (such as an unequal number of zeros and ones) then its stego capacity will be lower, as only a small proportion of candidate ciphertexts would look random enough.

D. The power of parity

We mentioned systems that generate a number of candidate locations for a ciphertext bit and then filter out the locations where actually embedding a bit would have a significant effect on the statistics thought to be relevant (in the case of hiding in an image, this could mean avoiding places where the local variance in luminosity is either very low or very high).

Our selection channel approach led us to suggest a better way [3]. We use our one-time pad (or keystream generator) to select not one pixel but a set of them, and embed the ciphertext bit as their parity. This way, the information can be hidden by changing whichever of the pixels can be changed least obtrusively.

From the information theoretic point of view, if each bit of the covertext is “1” with probability 0.6, then the probability that a bit pair will have parity 1 is 0.52; if we move to triples, the parity is 1 with probability 0.504, and so on. Thus by encoding each embedded bit as the parity of k bits of stegotext, we can reduce the effect that the embedding process has on the statistics of the stegotext below any arbitrary threshold; and as the improvement is geometric, we will not in practice have to increase k very much.

There is an interesting tradeoff: the more bits in the selection channel (i.e., the greater the value of k), the more bits we can hide in the covertext. In practice our selection channel will be a cryptographic pseudorandom number generator, and we can draw from it as many bits as we like.

There are still limits. For example, suppose that there is an allowed set of cover texts M (we might be using the cover of a news agency; we have to report a reasonably truthful version of events, and transmit photographs — perhaps slightly doctored—of events that actually took place). Suppose also that there is an allowed set of encodings E , and that each hidden bit is embedded by a choice of an encoding rule (such as a parity check in the method described above). Then the covert capacity will be at most $H(E) - H(M)$. But this gives us an upper bound only; it does not give us useful information on how much information may safely be hidden.

E. Equivalence classes

Suppose Alice uses a keyed cryptographic hash function to derive one bit from each sentence of a document. She may even have a macro in her word processor that checks every sentence as she finishes typing it and beeps if the output of the cryptographic hash function is not equal to the next bit of the message she wishes to embed. This alarm will go off about every other sentence, which she can then rewrite.

If she just uses standard synonym pairs such as [is able \leftrightarrow can], then clearly she must not alter their statistics to the point that Willie can detect the change. It is even an open question whether a computer can alter a natural language text in a way that is undetectable to a human [49]—that is, embed a ciphertext using the technique described above—and the problem is commended to the research community as the “Stego Turing Test.” Conversely, writing a program to scan for human inserted steganography might be rather hard. Recent work on natural language based stego is described in [54].

The use of synonyms to encode embedded messages is a special case of using equivalence classes of messages; these can also arise naturally in other applications. For example, when making a map from a larger scale map, many arbitrary decisions have to be taken about which details to incorporate, especially with features such as coastlines that are to some extent fractal [34]. Also, when software is written, it contains “birthmarks” such as the order in which registers are pushed and popped, and these were used by

IBM in litigation against software pirates who had copied their PC-AT ROM [22].

Equivalence classes of messages are tied up with compression. If covertext C_1 has a meaning or effect that is equivalent to that of covertext C_2 , then a compression algorithm need only select one representative from this equivalence class. However, if $C_1 \neq C_2$, then this choice throws away information, and the compression is lossy. Again, we get a bound on the stego channel capacity: it is the difference between lossy and lossless compression. Once more though, this is an upper bound rather than a safety bound, and is not much help against a powerful opponent.

It must be said that not all steganographic techniques involve equivalence classes. It is possible to create a series of images each of which differs only imperceptibly from the next, but such that the starting and final images are clearly different. This is relevant to the case where the warden is allowed to insert only so much distortion into messages; beyond a certain limit he might be held, in the absence of any hard evidence of covert activity by a prisoner, to have violated that prisoner’s human rights.

A purist might conclude that the only circumstance in which Alice can be certain that Willie cannot detect her messages is when she uses a true subliminal channel (see [7], and papers in this volume). However, other interesting things can be said about steganography.

V. ACTIVE AND PASSIVE WARDENS

We pointed out above that while an attack on a classical steganographic system consisted of correctly detecting the presence of embedded matter, an attack on a copyright marking scheme consists of rendering the mark useless.

There is a critical distinction between passive wardens, who monitor traffic and signal to some process outside the system if unauthorised traffic is detected, and active wardens who try to remove all possible covert messages from traffic that passes through their hands. In classical systems, the wardens could be either active or passive; while in marking systems, we are usually concerned with active wardens such as software pirates.

Consider the marking of executable code. Software birthmarks, as mentioned above, have been used to prove the authorship of code in court. They were more or less “automatically” generated when system software was hand assembled, but they must be produced more deliberately now that most code is compiled. One technique is to deliberately mangle the object code: the automatic, random replacement of code fragments with equivalent ones is used by Intel to customise security code [9].

One can imagine a contest between software authors and pirates to see who can mangle code most thoroughly without affecting its performance too much. If the author has the better mangler, then some of the information he adds will be left untouched by the pirate; but if the pirate’s code mangler is aware of all the equivalences exploited by the author’s, he may be able to block the stego channel completely. In general, if an active warden’s model of the communication is as good as the communicating par-

ties' model, and the covertext information separates cleanly from the usable redundancy, then he can replace the latter with noise.

In many other cases, the stego channel is highly bound up with the covertext. There have been measurements of the noise that can be added to a .gif file before the image quality is perceptibly degraded [23], and of the noise that can imperceptibly be added to digitised speech [16].

The point here is that if Alice can add an extra X% of noise without affecting the picture, then so can Willie; but where the noise and the signal do not separate cleanly, then both of them are limited to marking bits that are selected according to some rule or at random. Here, Alice can stop Willie finding out which X% carries the covert message by using a keystream to select the bits she will mark. In this case, all Willie will be able to do is to cut the bandwidth of the channel by adding his own random marks—a scenario that has been explored in the context of covert channels in operating systems [52].

This bandwidth limitation will also be effective against systems that embed each ciphertext bit as a parity check of a number of stegotext bits. When the warden is active, the more covertext bits we use in each parity check, the more easily he will be able to inject noise into our stegotext.

VI. PUBLIC KEY STEGANOGRAPHY

Until recently, it was generally assumed that, in the presence of a capable motivated opponent, steganography required the pre-existence of a shared secret so that the two communicating parties could decide which bits to tweak.

A. With a passive warden

In an early (workshop) version of this paper [3], we showed that public-key steganography is possible in the presence of a passive warden. Given a covertext in which any ciphertext at all can be embedded, then there will usually be a certain rate at which its bits can be tweaked without the warden noticing. So suppose that Alice can modify at least one out of every k bits of the covertext. This means that Willie cannot distinguish the parity of each successive block of k bits from random noise, and it follows that Alice can encode an arbitrary pseudorandom string in these parities.

This pseudorandom material lies in plain sight; anyone can read it. So Willie cannot tell the difference between stegotext and pure covertext by randomness tests; a suitable parity check function will extract pseudorandom-looking data from any covertext in which information could have been embedded at all.

Now suppose that Alice and Bob did not have the opportunity to agree a secret key before they were imprisoned, but that Bob has a public key that is known to Alice. She can take her covert message, encrypt it under his public key, and embed it as the parity of successive blocks. Each possible recipient will then simply try to decrypt every message he sees, and Bob alone will be successful. In practice, the value encrypted under a public key could be a control block consisting of a session key plus some padding, and

the session key would drive a conventional steganographic scheme.

Normal public key cryptography means that users can communicate confidentially in the absence of previously shared secrets; our construction of public key steganography shows that they can also communicate covertly (if this is at all possible for people who already share a secret).

B. With an active warden

The open question left in [3] was whether public key steganography was possible in the presence of a warden who is active rather than passive. The original construction fails in this case, as Willie can also tweak one bit in every k ; he could even set the parity of each successive block to zero. We will now indicate how concealed public key communication may still be possible in the presence of an active warden provided that the model of the Prisoners' Problem is changed slightly.

We will assume that the stegomessages Alice sends to Bob will be sent to other recipients too, such as a mailing list or usenet newsgroup. We will also assume that Willie and Alice are each able to tweak at most one bit in k of the content (as above, Willie might infringe Bob's rights if he distorts the communication channel to the point that it becomes unusable). Alice can choose a short one-time key that selects some permutation of the covertext bits, and she hides a message as the parity of successive k -tuples of bits in this permuted sequence. Willie, suspecting that this method may be in use, alters 1 in k of the stegotext bits; this is the best he can do since he does not know what one-time key Alice used. This corrupts most of the bits in Alice's message, but not all of them; asymptotically, about $1/e$ k -tuples will be unaffected, and so there will be a positive residual channel capacity. Given a suitable error correcting code, Alice can still send a message encrypted using Bob's public key.

Once Bob has received the message, Alice broadcasts her short one-time key. Bob now applies it to all the messages he has in store; lo and behold, one of them produces a bit string that he can decrypt using his private key. Willie can also now tell that one of the messages he forwarded from Alice to Bob contained suspicious content, namely a random looking string with an error correction code attached, that was most likely an instance of the protocol described here. However, he cannot tell that the message was directed specifically to Bob, as he does not possess Bob's private key.

We have changed our model slightly, in that we now assume that Alice can send short pseudorandom messages with integrity. Otherwise, when Willie sees the message with the one-time key and realises its significance, he would corrupt it or refuse to forward it. However, there might be circumstances in which he is unable to do this. For example, Alice might be at liberty while Bob is in jail; and Willie might be able to censor Alice's usenet postings via the prison's news server, but not permitted to censor comments that she makes from time to time on radio programmes.

Our idea led to a suggestion of another approach, in which a slightly different change is made in the assumptions of the Prisoners' Problem—namely that tamper-resistant embedding is possible. A message encrypted under Bob's public key (or, alternatively, Alice's short one-time key) might be embedded as a high-level description of the cover object in such a way that it could not easily be removed [15]. For example, one might encode the message as the locations mentioned in a short story, together with the order of their appearance. It is clearly possible for an author to so entwine known features of towns and countries into a narrative, that any attempt to change them would require a complete rework of the plot.

Both of these methods may appear more contrived than practical, and they serve more as existence proofs than as practical engineering proposals. They also serve to emphasise that very small changes in our starting assumptions can have a significant effect on the conditions under which we can hide information.

VII. CONCLUSIONS

We have explored the limits of steganographic theory and practice. We started off by outlining a number of techniques both ancient and modern, together with attacks on them (some new); we then discussed a number of possible approaches to a theory of the subject. We pointed out the difficulties that stand in the way of a theory of "perfect covertness" with the same power as Shannon's theory of perfect secrecy. But considerations of entropy give us some quantitative leverage and the "selection channel"—the bandwidth of the stego key—led us to suggest embedding information in parity checks rather than in the data directly. This approach gives improved efficiency, and also allows us to do public key steganography. Finally, we have shown that public key steganography may sometimes be possible in the presence of an active warden.

Acknowledgements: Some of the ideas presented here were clarified by discussion with David Wheeler, John Daugman, Roger Needham, Gus Simmons, Markus Kuhn, Peter Rayner, David Aucsmith, John Kelsey, Ian Jackson, Mike Roe, Mark Lomas, Stewart Lee, Peter Wayner, Matt Blaze and Scott Craver. The second author is grateful to Intel Corporation for financial support under the grant "Robustness of Information Hiding Systems."

REFERENCES

- [1] "Auditory masking and MPEG-1 audio compression," E Ambikairajah, AG Davis, WTK Wong, *IEE Electronics & Communication Engineering Journal* v 9 no 4 (Aug 97) pp 165-175
- [2] "Liability and Computer Security: Nine Principles", RJ Anderson, in *Computer Security—ESORICS 94*, Springer LNCS v 875 pp 231-245
- [3] "Stretching the Limits of Steganography", RJ Anderson, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 39-48
- [4] "The Eternity Service", in *Proceedings of Pragocrypt 96* (GC UCMP, ISBN 80-01-01502-5) pp 242-252
- [5] "Tamper Resistance—a Cautionary Note", RJ Anderson, MG Kuhn, in *Proceedings of the Second Usenix Workshop on Electronic Commerce* (Nov 96) pp 1-11
- [6] "Chameleon—A New Kind of Stream Cipher", R Anderson, C Manifavas, to appear in *Proceedings of the 4th Workshop on Fast Software Encryption* (1997)
- [7] "The Newton Channel", RJ Anderson, S Vaudenay, B Preneel, K Nyberg, *this volume*
- [8] www.anonymizer.com
- [9] "Tamper Resistant Software: An Implementation", D Aucsmith, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 317-333
- [10] "Techniques for Data Hiding", W Bender, D Gruhl, N Morimoto, A Lu, *IBM Systems Journal* v 35 no 3-4 (96) pp 313-336
- [11] "Watermarking Digital Images for Copyright Protection", FM Boland, JJK Ó Ruanaidh, C Dautzenberg, *Proceedings, IEE International Conference on Image Processing and its Applications, Edinburgh 1995*
- [12] "Digital Watermarks for Audio Signals," L Boney, AH Tewfik, KN Hamdy, in *IEEE International Conference on Multimedia Computing and Systems*, June 17-23, 1996 Hiroshima, Japan; pp 473-480
- [13] "A Secure, Robust Watermark for Multimedia", IJ Cox, J Kilian, T Leighton, T Shamoon, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 183-206
- [14] R Cox, presentation at the 'Access All Areas' conference, London, UK, May 7 1997
- [15] "On Public-key Steganography in the Presence of an Active Warden", S Craver, *IBM Research Report RC 20931*, July 23, 1997
- [16] "Computer Based Steganography", E Franz, A Jerichow, S Möller, A Pfitzmann, I Stierand, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 7-21
- [17] "Hiding Routing Information", DM Goldschlag, MG Reed, PF Syverson, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 137-150
- [18] "Echo Hiding", D Gruhl, A Lu, W Bender, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 295-315
- [19] 'Copyright theft', J Gurnsey, Aslib Gower, 1995
- [20] "A voluntary international numbering system—the latest WIPO proposals", R Hart, *Computer Law and Security Report* v 11 no 3 (May-June 95) pp 127-129
- [21] 'Speech Synthesis and Recognition—Aspects of Information Technology', JN Holmes, Chapman & Hall, 1993
- [22] Talk on software birthmarks, counsel for IBM Corporation, BCS Technology of Software Protection Special Interest Group, London 1985
- [23] 'Steganography in Digital Images', G Jagpal, Thesis, Cambridge University Computer Laboratory, May 1995
- [24] 'The Codebreakers', D Kahn, Macmillan 1967
- [25] 'La Cryptographie Militaire', A Kerckhoffs, *Journal des Sciences Militaires*, 9th series, IX (Jan 1883) pp 5-38; (Feb 1883) pp 161-191
- [26] "Towards Robust and Hidden Image Copyright Labeling", E Koch, J Zhao, *Proceedings of 1995 IEEE Workshop on Non-linear Signal and Image Processing* (Neos Marmaras, Halkidiki, Greece, June 20-22, 1995)
- [27] "Phreaking recognised by Directorate General of France Telecom", HM Kriz, in *users/Chaos Digest 1.03 (Jan 93)*
- [28] "A Cautionary Note on Image Downgrading", C Kurak, J McHugh, *Computer Security Applications Conference*, (IEEE, 1992) pp 153-159
- [29] 'Codes, Keys and Conflicts: Issues in U.S. Crypto Policy', S Landau, S Kent, C Brooks, S Charney, D Denning, W Diffie, A Lauck, D Miller, P Neumann, D Sobel, Report of a Special Panel of the ACM U.S. Public Policy Committee, June 1994
- [30] "Copy Protection for Multimedia Data based on Labeling Techniques", GC Langelaar, JCA van der Lubbe, J Biemond, 17th Symposium on Information Theory in the Benelux, Enschede, The Netherlands, May 1996
- [31] "Electronic Document Distribution", NF Maxemchuk, *AT & T Technical Journal* v 73 no 5 (Sep/Oct 94) pp 73-80
- [32] 'An Introduction to the Psychology of Hearing', BCJ Moore, Academic Press 1989
- [33] "Covert Channels—Here to Stay?", IS Moskowitz, MH Kang, *Compass 94* pp 235-243
- [34] RM Needham, *private conversation*, December 1995
- [35] 'Voice and Speech Processing', T Parson, McGraw-Hill, 1986
- [36] "Information Hiding Terminology", B Pfitzmann, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 1-11

- tion Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 347–350
- [37] “Trials of Traced Traitors”, B Pfitzmann, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 49–64
- [38] “A method for signature casting on digital images”, I Pitas, in *International Conference on Image Processing* v 3 (Sep 96) pp 215–218
- [39] “Crowds: Anonymity for web transactions”, MK Reiter, AD Rubin, *DIMACS Technical Report 97-15* (April 1997)
- [40] ‘*Meteor Burst Communications: Theory and Practice*’, DL Schilling, Wiley 93
- [41] ‘*Applied Cryptography—Protocols, Algorithms and Source Code in C*’ B Schneier (second edition), Wiley 1995
- [42] “A Mathematical Theory of Communication”, CE Shannon, in *Bell Systems Technical Journal* v 27 (1948) pp 379–423, 623–656
- [43] “Communication theory of secrecy systems”, CE Shannon, in *Bell Systems Technical Journal* v 28 (1949) pp 656–715
- [44] “The Prisoners’ Problem and the Subliminal Channel”, GJ Simmons, in *Proceedings of CRYPTO ’83*, Plenum Press (1984) pp 51–67
- [45] “How to Insure that Data Acquired to Verify Treaty Compliance are Trustworthy”, GJ Simmons, *Proceedings of the IEEE* v 76 (1984) p 5
- [46] “A survey of information authentication”, GJ Simmons, in *Contemporary Cryptology—the Science of Information Integrity*, IEEE Press 1992, pp 379–419
- [47] “The History of Subliminal Channels”, GJ Simmons, *this volume*
- [48] ‘*High Quality De-interlacing of Television Images*’, N van Someren, PhD Thesis, University of Cambridge, September 1994
- [49] K Spärck Jones, *private communication*, August 1995
- [50] “Police to shut out snoopers”, *Sunday Times* (13 July 1997) p 3.13
- [51] www.surety.com
- [52] “Modelling a Fuzzy Time System”, JT Trostle, *Proc. IEEE Symposium in Security and Privacy 93* pp 82 - 89
- [53] “A Digital Watermark”, RG van Schyndel, AZ Tirkel, CF Osborne, in *International Conference on Image Processing*, (IEEE, 1994) v 2 pp 86–90
- [54] ‘*Disappearing Cryptography—Being and Nothing on the Net*’, P Wayner, AP Professional (1996)
- [55] “Fighting Mobile Phone Fraud—Who is Winning?”, K Wong, in *Datenschutz und Datensicherung*, pp 349–355, June 1995



Ross J. Anderson received his BA, MA and PhD degrees from the University of Cambridge, UK, where he is a university lecturer at the Computer Laboratory. He teaches and directs research in computer security and software engineering, and was the program chair of the first international workshop on information hiding, held at Cambridge in May–June 1996. He is a Fellow of the RSA and the IMA, a Member of the IEE, and a Chartered Engineer; he is also the editor of ‘Computer and Communications Security Reviews’. His research interests centre on security engineering: he has published extensively on how computer security systems fail and what can be done to make them more robust.



Fabien A.P. Petitcolas graduated from the École Centrale, Lyon, France and received a Diploma in Computer Science from the University of Cambridge, UK, where he is currently a research student at the Computer Laboratory. His research topic is the robustness of information hiding systems.