

ON THE LIMITS OF STEGANOGRAPHY

R J Anderson

Rapporteur: Ian Welch

On The Limits of Steganography

Ross Anderson, Fabien Petitcolas

Cambridge University Computer Laboratory
 Pembroke Street, Cambridge CB2 3QG, UK
 Email (rja14,fapp2)@cl.cam.ac.uk

Abstract. In this paper, we seek to 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.

1 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, making a secret message visible. [23].

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. During the first and second world wars, postal censors deleted lovers' X's, shifted watch hands, and replaced items such as loose stamps and blank paper. They also rephrased telegrams; in one case, a censor changed 'father is dead' to 'father is deceased', which elicited the reply 'is father dead or deceased?'

The study of this subject in the scientific literature may be traced to Simmons, who in 1983 formulated it as the "Prisoners' Problem" [42]. 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 of documents, 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 [30].

Simmons' formulation of the Prisoners' Problem was itself an instance of information hiding. It was a ruse to get the academic community to pay some 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 information past the other country's monitoring facilities [43, 45].

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 [48].

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; the main practical problem facing law enforcement and intelligence agencies is acknowledged to be real-time 'traffic selection' — deciding in real time which of the huge mass of calls to intercept [28]. Criminal techniques 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 [26]; while in the UK case, a drug dealer physically tapped into a neighbour's phone [13]. All these techniques also involve theft of service, from either the phone company or 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 [52].

Military organisations also use unobtrusive communications. Their preferred mechanisms include spread spectrum and meteor scatter radio [38], which can give various possible combinations of resistance to detection, direction finding and jamming and are important for battlefield communications where the radio operators are at risk of being located and attacked. Similar facilities are available on the Internet, where 'anonymous remailers' can be used to hide the origin of an electronic mail message, and are being developed for other protocols such as ftp and http [7, 16, 37].

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.

2 What is Steganography?

Classical steganography concerns itself with techniques for 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 another key related to it) 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 referred to as 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 [35]).

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 [15].

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 videomail, 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 would 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 to prove ownership later. One may think of a watermark as one or more copyright marks that are hidden in the content, and that may be displayed in court to prove a case of piracy.

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 an 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 [5].

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.

But 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 [36] or involving a public timestamping service in the marking process.

3 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 Kerkhoffs in 1883 [24], and has been borne out by long and hard experience since [23]. It should be particularly obvious in applications where the protection mechanism is to provide evidence [1].

However, most promoters of copyright marking and steganographic 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.

3.1 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 [51]. Information embedded in this way may be invisible to the human eye [27] 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 [39] 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 [15].

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). A bit can be embedded in a pixel that passes this test by some rule such as setting the pixel's least significant bit to the parity of the surrounding pixels to embed a '1'.

Such schemes can be tricky to implement (for example, one must ensure that different blocks of nine neighbouring pixels do not overlap and interfere with each other). They can also be destroyed in a number of ways by an opponent who can modify the stegoimage. 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 differential luminance of a large number of pseudorandomly chosen pixel pairs [9]; similar techniques can be used for digital audio.

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 to work slightly harder 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 amplitude, which is tiny compared to the precision with which the original sounds were in most cases generated.

This technique is very crude compared with what one could do using more sophisticated resampling and filtering algorithms. Yet we still found 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.

3.2 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 threat 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) [22], 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 [25, 29] 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 [12]; 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 [49] in conjunction with the marking system, so that the priority of the genuine mark can be independently established). Other schemes of this kind include, for example, [10] and [25].

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 [11].

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 [20, 34]. 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 a wideband noise, only a small band centred about the tone contributes to the masking effect [31]. 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, 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 [11]. This gives an example of a mark whose existence is detectable by statistical tests and yet still be undetectable to humans; whether it can be damaged beyond later recognition without introducing perceptible distortion is an interesting question.

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 [17]: this 'echo hiding' technique manipulates the cepstral transform of an audio signal in order to embed an echo. This echo might have a delay of 0.5 ms to signal a '0' and 1.0 ms to signal a '1'; and it has been found experimentally that data can be hidden at up to 15 bits per second in modern music.

3.3 A general model

So the general model of steganography is that Alice embeds information by tweaking some bits of some transform of the coverttext. This transform enables

her to get at one or more bits which are redundant in the sense that when she tweaks some subset of them, this cannot be detected easily or at all by an opponent who does not know which subset to look at.

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. The usual result is that 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 able and willing to do even more work. Thus the apparent redundancy which we exploit will fall within his ability to model. This may be especially the case 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 music recording in 1997, 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" [18]. Such experience is emphasised by the recent success of criminals in cloning the smartcards used to control access to satellite TV systems [4].

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 [41]. Simmons provided us with a comparable theory of authentication, that has been applied in nuclear weapons command and control [44]. 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.

4 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 coverttext, 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 coverttext's intrinsic noise. Is there any way this intuitive result could be proved?

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) [23]. 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 coverttext 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 coverttext 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.

4.4 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 [2]. 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 the coverttext is '1' with probability 0.6, and we encode in bit pairs, 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. As the improvement is geometric, we can get the discrepancy as low as we like.

There is an interesting tradeoff: the more bits in the selection channel, the more bits we can hide in the coverttext. 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 . 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.

4.5 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 Willy 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 [47] — 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.

This is an example of equivalence classes of messages; these can often arise naturally. 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 [33]. 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 [21].

Equivalence classes of messages are tied up with compression. If coverttext C_1 means the same as coverttext 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 so 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 [6], and papers in this volume). However, other interesting things can be said about steganography.

5 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 so that pirates could be prosecuted. They were more or less 'automatically' generated with hand assembled system software, but must be produced more deliberately now that most code is compiled. One technique is to deliberately mangle the object code in some way: the automatic, random replacement of code fragments with equivalent ones is used by Intel to customise security code [8].

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 parties', 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 [22], and of the noise that can imperceptibly be added to digitised speech [15].

The point here is that if Alice can add an extra $X\%$ of noise without affecting the picture, then so can Willie; but she can stop him finding out which $X\%$ carries the covert message by using a keystream to select which changes to make. In this case, all Willie will be able to do is to cut the bandwidth of the channel — a scenario that has been explored in the context of covert channels in operating systems [50].

This bandwidth limitation will also be effective against systems that embed each ciphertext bit as a parity check of a number of covertext 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 covertext.

6 Public Key Steganography — Revisited

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.

In 1996, we showed that that public-key steganography is possible in the presence of a passive warden [2]. 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 she can encode an arbitrary pseudorandom string in these parities.

This pseudorandom material lies in plain sight; anyone can read it. So Willie cannot test the stegotext by applying randomness tests to it; a suitable parity check function will extract pseudorandom-looking data from any message in which covert information can be inserted 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).

The open question left in [2] 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, he might be permitted to censor Alice's usenet postings but not comments that she makes from time to time on radio programmes. (Of course, a spoken radio broadcast made in a studio that Alice did not control might have a very low stego capacity, hence the need for the bulk of the message to be sent via usenet.)

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 [14]. For example, one might encode the message as the names of the geographical 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. However, they make the point that very small changes in the starting assumptions can have a significant effect on the conditions under which we can hide information.

7 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 and Matt Blaze. The second author is grateful to Intel Corporation for financial support under the grant 'Robustness of Information Hiding Systems'.

References

1. "Liability and Computer Security: Nine Principles", RJ Anderson, in *Computer Security — ESORICS 94*, Springer LNCS v 875 pp 231-245
2. "Stretching the Limits of Steganography", RJ Anderson, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 39-48
3. "The Eternity Service", in *Proceedings of Pragocrypt 96* (GC UCMP, ISBN 80-01-01502-5) pp 242-252
4. "Tamper Resistance — a Cautionary Note", RJ Anderson, MG Kuhn, in *Proceedings of the Second Usenix Workshop on Electronic Commerce* (Nov 96) pp 1-11
5. "Chameleon — A New Kind of Stream Cipher", R Anderson, C Manifavas, to appear in *Proceedings of the 4th Workshop on Fast Software Encryption (1997)*
6. "The Newton Channel", RJ Anderson, S Vaudenay, B Preneel, K Nyberg, *this volume*
7. www.anonymizer.com
8. "Tamper Resistant Software: An Implementation", D Aucsmith, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 317-333
9. "Techniques for Data Hiding", W Bender, D Gruhl, N Morimoto, A Lu, *IBM Systems Journal* v 35 no 3-4 (96) pp 313-336
10. "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*
11. "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
12. "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
13. R Cox, talk to 'Access All Areas' conference, London, 5/7/97
14. S Craver, *private communication*
15. "Computer Based Steganography", E Franz, A Jerichow, S Moeller, A Pfitzmann, I Stierand, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 7-21
16. "Hiding Routing Information", DM Goldschlag, MG Reed, PF Syverson, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 137-150
17. "Echo Hiding", D Gruhl, A Lu, W Bender, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 295-315
18. 'Copyright theft', J Gurnsey, Aslib Gower, 1995
19. "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
20. 'Speech Synthesis and Recognition — Aspects of Information Technology', JN Holmes, Chapman & Hall, 1993
21. Talk on software birthmarks, counsel for IBM Corporation, BCS Technology of Software Protection Special Interest Group, London 1985
22. 'Steganography in Digital Images', G Jagpal, Thesis, Cambridge University Computer Laboratory, May 1995
23. 'The Codebreakers', D Kahn, Macmillan 1967
24. "La Cryptographie Militaire", A Kerkhoffs, *Journal des Sciences Militaires*, 9th series, IX (Jan 1883) pp 5-38; (Feb 1883) pp 161-191

25. "Towards Robust and Hidden Image Copyright Labeling", E Koch, J Zhao, *Proceedings of 1995 IEEE Workshop on Nonlinear Signal and Image Processing* (Neos Marmaras, Halkidiki, Greece, June 20-22, 1995)
26. "Phreaking recognised by Directorate General of France Telecom", HM Kriz, in *Chaos Digest 1.03 (Jan 93)*
27. "A Cautionary Note on Image Downgrading", C Kurak, J McHugh, *Computer Security Applications Conference*, (IEEE, 1992) pp 153-159
28. '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
29. "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
30. "Electronic Document Distribution", NF Maxemchuk, *AT & T Technical Journal* v 73 no 5 (Sep/Oct 94) pp 73-80
31. 'An Introduction to the Psychology of Hearing', BCJ Moore, Academic Press 1989
32. "Covert Channels — Here to Stay?", IS Moskowitz, MH Kang, *Compass 94* pp 235-243
33. RM Needham, *private conversation*, December 1995
34. *Voice and Speech Processing*, T Parson, McGraw-Hill, 1986
35. "Information Hiding Terminology", B Pfitzmann, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 347-350
36. "Trials of Traced Traitors", B Pfitzmann, in *Information Hiding*, Springer Lecture Notes in Computer Science v 1174 (1996) pp 49-64
37. "Crowds: Anonymity for web transactions", MK Reiter, AD Rubin, *DIMACS Technical Report 97-15* (April 1997)
38. 'Meteor Burst Communications: Theory and Practice', DL Schilling, Wiley 93
39. 'Applied Cryptography — Protocols, Algorithms and Source Code in C' B Schneier (second edition), Wiley 1995
40. "A Mathematical Theory of Communication", CE Shannon, in *Bell Systems Technical Journal* v 27 (1948) pp 379-423, 623-656
41. "Communication theory of secrecy systems", CE Shannon, in *Bell Systems Technical Journal* v 28 (1949) pp 656-715
42. "The Prisoners' Problem and the Subliminal Channel", GJ Simmons, in *Proceedings of CRYPTO '83*, Plenum Press (1984) pp 51-67
43. "How to Insure that Data Acquired to Verify Treaty Compliance are Trustworthy", GJ Simmons, *Proceedings of the IEEE* v 76 (1984) p 5
44. "A survey of information authentication", GJ Simmons, in *Contemporary Cryptology — the Science of information Integrity*, IEEE Press 1992, pp 379-419
45. "The History of Subliminal Channels", GJ Simmons, *this volume*
46. 'High Quality De-interlacing of Television Images', N van Someren, PhD Thesis, University of Cambridge, September 1994
47. K Spärck Jones, *private communication*, August 1995
48. "Police to shut out snoopers", *Sunday Times* (13 July 1997) p 3.13
49. www.surety.com
50. "Modelling a Fuzzy Time System", JT Trostle, *Proc. IEEE Symposium in Security and Privacy 93* pp 82 - 89
51. "A Digital Watermark", RG van Schyndel, AZ Tirkel, CF Osborne, in *International Conference on Image Processing*, (IEEE, 1994) v 2 pp 86-90
52. "Fighting Mobile Phone Fraud — Who is Winning?", K Wong, in *Datenschutz und Datensicherung (6/95)* pp 349-355

DISCUSSION

Rapporteur : Ian Welch

Lecture One

A participant asked if the two images shown (one with another image contained within it) are different (they look different on the OHP). Dr Anderson suggested that the difference was due to the printing of the image. On the screen the differences were imperceptible.

Professor Lobelle asked if one can conclude it is always possible to attack schemes designed to hide data within audio. If so then this is a very negative conclusion, is it ever possible to prove it impossible to break. Dr Anderson replied that it is impossible to prove that you cannot break such a scheme but it is possible to offer the user a trade-off between distortion and bandwidth. There is no resistance against attack if you understand the key. Professor Lobelle suggested that perhaps a solution was to use a one-time key. Dr Anderson replied that even if a one-time key was used it would still be possible to carry out a successful attack.

