

Loud and Clear: Human-Verifiable Authentication Based on Audio

Michael T. Goodrich, Michael Sirivianos, John Solis, Gene Tsudik, and Ersin Uzun
Department of Computer Science
University of California, Irvine
{goodrich,msirivia,jsolis,gts,euzun}@ics.uci.edu

Abstract

*Secure pairing of electronic devices that lack any previous association is a challenging problem which has been considered in many contexts and in various flavors. In this paper, we investigate the use of the audio channel for human-assisted authentication of previously un-associated devices. We develop and evaluate a system we call **Loud-and-Clear (L&C)** which places very little demand on the human user. L&C involves the use of a text-to-speech (TTS) engine for vocalizing a robust-sounding and syntactically-correct (English-like) sentence derived from the hash of a device’s public key. By coupling vocalization on one device with the display of the same information on another device, we demonstrate that L&C is suitable for secure device pairing (e.g., key exchange) and similar tasks. We also describe several common use cases, provide some performance data for our prototype implementation and discuss the security properties of L&C.*

Index Terms: *Human-assisted authentication, Man-in-the-middle attacks, Audio, Text-to-speech, Public key, Personal device, Wireless networks.*

1 Introduction

The proliferation of many types of inexpensive personal devices, such as PDAs, cellphones, smart watches, and MP3 players, has been accompanied by the need to secure these devices and their communication with the “outside world.” Common applications involve securely connecting one’s personal device to an unfamiliar printer, fax machine, wireless projector, network access point or another personal device. Of course, establishing such secure connections is straightforward if there exists a pervasive global security infrastructure, such as a public key infrastructure (PKI) or a trusted online third party (TTP). However, in many (even most) application scenarios involving heterogeneous personal devices, neither a PKI nor a TTP can be assumed. Thus we are faced with the with the problem of *peer device authentication* or *secure device pairing*.

While the general problem of secure pairing of devices with no prior context is difficult and remains partially unsolved, there has been progress in scenarios involving personal devices. Precisely because these devices are *personal*, the presence of the human user (owner) is assured and techniques have been proposed to engage the human user in the process of establishing secure communication. Therefore, human assistant authentication represents a very timely, important and popular research topic.

In this paper, we focus on application settings where devices are physically present and near but their communication channel is not visible or ascertainable by their owners or users. The most obvious example is whenever nearby devices communicate over a wireless channel, e.g., 802.11a/b/g, Bluetooth, or Infrared. Such channels offer no physical evidence of direct connection between devices.¹ To address ad hoc secure pairing of devices over such communication channels, we develop a system called **Loud-and-Clear (L&C)** which uses the audio channel to attain human-assisted (but not burdensome) device authentication.

Organization. This paper is organized as follows: The next section overviews related work. It is followed by the discussion of our motivation and the summary of our contributions in Section 3. Then, Section 4 describes the key elements of

¹Although we use wireless communication as a running motivation throughout this paper, the difficulty of peer device authentication is not confined to wireless links. For example, if the communication between the two devices is via wired Ethernet, a similar problem arises. Only a fully visible point-to-point physical connection would alleviate the peer authentication problem we study in this paper.

L&C. Sections 5 and 6 discuss related topics of unidirectional authentication and presence confirmation, respectively. L&C prototype implementation and its performance are discussed in Section 7 and security issues are considered in Section 8.

2 Related Work

There has been a considerable amount of prior work in the general area of secure device pairing. As mentioned earlier, PKI-oriented solutions involve rigid hierarchies and require the existence of trusted off-line Certification Authorities (CAs) [1]. Unstructured certification techniques such as PGP [2] assume a certain small degree of separation among certified entities (i.e., a web of trust). An alternative is an online trusted third party (TTP), such as Kerberos [3, 4]. However, such solutions require constant presence and availability of an online TTP which is not realistic in our envisaged scenarios. For this reason, the remainder of this section focuses on closely related prior research.

The well-known Diffie-Hellman key exchange protocol [5] allows two entities, with no prior secrets or secure association, to agree on a common secret key. However, precisely because no prior secrets are assumed, man-in-the-middle (MITM) attacks are possible [6, 7, 8]. A number of enhancements to Diffie-Hellman protocol have been developed. Bellare and Merritt proposed the encrypted key exchange (EKE) protocol [9] to prevent MITM attacks. However, it requires both parties to possess a secret password *a priori*. Although many EKE refinements have been proposed [10, 11, 12], all involve a pre-shared secret password. This is clearly inapplicable in our targeted environment.

Secondary channels offer another means to defend against MITM attacks. A secondary channel can be used to verify that the keys computed at both devices are identical. A number of efforts have been made to involve a human user in the secondary channel in order to manually verify/compare keys (or hashes thereof). Unfortunately, since a numerically-represent hash can be long and cumbersome, comparing hashes is an arduous process. To ease the burden, visual metaphors to represent hash values have been devised, including the use of fractal snowflakes [13, 14], random art [15] and flag representation [16]. These representations allow for human users to do rough similarity checking of strings, but they do not allow for fine-grained comparisons, as required in many security applications.

One notable recent result is the *Seeing-is-Believing* (SiB) system proposed by McCune, et al. [17]. SiB takes advantage of the visual channel – since a human user is assumed capable of visually identifying the target device – to provide human-assisted authentication. This is done by requiring the human user to take a picture (with a personal camera-equipped mobile phone) of the 2D barcode affixed to, or displayed by, a target device and having the phone interpret the barcode and extract the cryptographic material identifying the device’s public key. Meanwhile, the target device is supposed to communicate to the user’s phone the (presumably) same cryptographic material via some wireless channel, e.g., Bluetooth. If the two versions of the cryptographic material match, the user’s phone concludes that the target device’s public key is authentic.

SiB provides a reasonable level of security, commensurate with what is realistic under the circumstances. Circumventing SiB requires the adversary to: (1) either hack into the target device and cause it to display wrong barcodes or physically plaster fake barcodes on the device, and, (2) mount a man-in-the-middle attack on the wireless channel. SiB is also quite practical particularly because it places very little burden upon the human user: visual identification of the target device and taking a picture of the barcode.

Another research direction similar to the one taken in this paper is the textual representation of cryptographic strings. For example, the S/KEY One-Time-Password [18] system represents an MD5 hash as a series of six short (up to four-letter) words. S/KEY was designed for the automatic generation of pass-phrases, and the pass-phrases are not necessarily auditorially robust nor syntactically-correct. That is, S/KEY pass-phrases were not meant to be spoken, and there are some similar-sounding words in the S/KEY word list. The word encoding for PGPfone [19], on the other hand, uses an auditorially robust word list. Still, it does not produce syntactically-correct string encodings, which makes it harder for human users to parse PGPfone phrases or even to read along with a spoken version of a PGPfone phrase. Our L&C system produces phrases that are syntactically-correct and auditorially robust; hence, they are more suitable for human verification of phrase similarity.

Also of interest are schemes for constructing syntactically-correct pass-phrases from random passwords, e.g., see [20]. Such techniques essentially allow users to remember long passwords using the first letters of the words in an easy-to-remember pass-phrase. While syntactically-correct, they are not necessarily auditorially robust.

For peer device authentication, Stajano and Anderson proposed a method for establishing keys by means of a link created through physical contact [21]. However, due to the diversity in personal devices, it is impractical to expect all devices to have suitable physical connection interfaces. Likewise, it may also be infeasible to lug around connection interfaces and or interface converters for various devices. Balfanz, et al. [22] extended this approach to location-limited channels through the use of short-range wireless Infrared communication. Capkun et al. [23] proposed a further extension to allow two previously unassociated devices to establish a key utilizing one-hop transitive trust.

Manual peer device authentication has also been studied in [24, 25, 26] and [27]. The common element of these solutions is that the user is required to compare short numerical check values, which are generated by hashing or taking the MAC of the authentication object. Their limitation is that sufficient security dictates that the check values need to be relatively lengthy (substantially more than the 4-digit hex vector [24] suggests) rendering their comparison a cumbersome and error-prone task for humans. [25] improves on [24] by shortening the required length of the check value, yet the user is burdened with the task of typing a short 2-4 hex digit key using the non-user-friendly input interface of a personal device. [27] also proposes a solution to reduce the length requirement for the check values. However, it requires a temporary shared secret between the two peers, thus it is not applicable for devices with no prior association. More recently, Cagalj et al [28] tackled the problem of user-friendly mutual authentication with three commitment-based [29] techniques. The first requires the user to compare approximately half the number of bits of the previous solutions. The other two techniques rely on the specifics of the radio channel; they use distance bounding and integrity codes.

3 Overview and Motivation

In this paper, we investigate the use of the audio channel for human-assisted authentication of unfamiliar devices. We develop and evaluate a system called **Loud-and-Clear (L&C)**, which, like Seeing-is-Believing (SiB), places relatively little demand on the human user. L&C uses spoken natural language for human-assisted authentication; hence, it is suitable for secure device pairing (e.g., key exchange) and similar tasks, such as secure configuration.

The motivation for our work is three-fold:

1. While some mobile phones include a built-in photo (and even video) camera, many other types of personal devices are not similarly equipped. For example, a camera is not standard equipment on most PDAs (e.g., Treos, Blackberries, IPAQs and PalmPilots). It is also not present on digital music players and smart watches. Furthermore, even for mobile phones, an on-board camera results in a certain price differential, as compared to a similar camera-less phone.
2. The use of cameras is appropriate for most people, except for those who are visually impaired (e.g., legally blind). One possibility is to print barcodes in Braille, ask the visually-impaired user to identify the target device's barcode by touch and then take the picture (with the camera phone) at very close range. Albeit, the associated burden would be higher than in plain SiB.
3. Barcodes and cameras can be used in many normal everyday settings, such as offices, hotels and airports. However, there are two important underlying assumptions: (1) ample ambient light, and, (2) sufficient proximity between the two devices. In other words, in the presence of light-inhibiting environmental factors, such as darkness, smoke or heavy fog, SiB would not be applicable.

By relying on the use of spoken natural language to provide human-verifiable secure communication, L&C alleviates all of the above shortcomings of SiB. Moreover, the L&C system encodes authentication strings using auditorially-robust, syntactically-correct "MadLib" phrases, which allows human users to easily verify the authentication strings between peer devices. To construct auditorially-robust text sequences, we produce a number of word lists of appropriate parts of speech, with the words in each list being as phonetically distant from each other as possible.

Indeed, we describe in this paper implementations of the L&C system applied to several peer device authentication scenarios. We also describe performance data for our implementation and security issues related to L&C. That is, we show that L&C provides an effective solution to many peer device authentication problems.

Nevertheless, the use of the audio channel for human-assisted authentication has its own drawbacks and limitations, which we readily admit from the outset:

1. Ambient noise is clearly an inhibiting factor for audio-based authentication. Whether comparing two audible sequences or comparing one such sequence to a displayed textual representation of the same sequence, noise makes authentication difficult. By the same token, the audio channel is not suitable for hearing-impaired human users.
2. As discussed later in this paper, L&C requires for at least one of the two devices to have a speaker (or an audio-out interface). The other device must either have a display or a speaker. While such interfaces are more common than cameras (as most personal devices are equipped with a speaker or a display, and often have both), we note that SiB does not require the use of a speaker or audio-out signal (it requires a camera/scanner and a display).
3. L&C places an alternate burden from SiB on the human user. In SiB, the user is asked to visually identify the target device and to take a picture of the device's barcode. In contrast, L&C requires the user to either compare two audio sequences or compare one audio sequence to a displayed textual representation of the same sequence.

It is apparent from the above that, owing to their respective advantages and limitations, SiB and L&C complement each other.

4 Main Elements of L&C

In this section, we describe the main elements of the Loud-and-Clear (L&C) system. We focus initially on the authentication of a target device to a personal device, assuming no prior context between the personal device and target device. We assume that in this, and most use scenarios, identification of the communicating devices is performed visually (or tactily) by the human user.

We consider the most plausible type of *authentication object*, which is the **Target device's public key**. The personal device receives the target device's public key over a wireless channel (e.g., Infrared, 802.11a/b/g, Bluetooth) and an audio signal is used as a means of verifying this public key or the one-way hash of this public key.

4.1 Requirements

The specifics of target device authentication in L&C depend upon several factors, such as the type of authentication objects, directionality, the number of human users, and the device equipment. The following basic requirements are common to all use cases:

- There is at least one human user present with a personal device.
- At least one device has an audio interface, e.g. a speaker or a audio out plug (though, as discussed below, for the sake of completeness, L&C also supports the case of both devices having displays but no audio).
- The two devices must be able to communicate via some multiple-access broadcast medium, e.g., 802.11a/b/g, Bluetooth, Infrared, or wired Ethernet. We make no assumptions regarding the security of this channel.

Figure 1 depicts some anticipated L&C use scenarios.

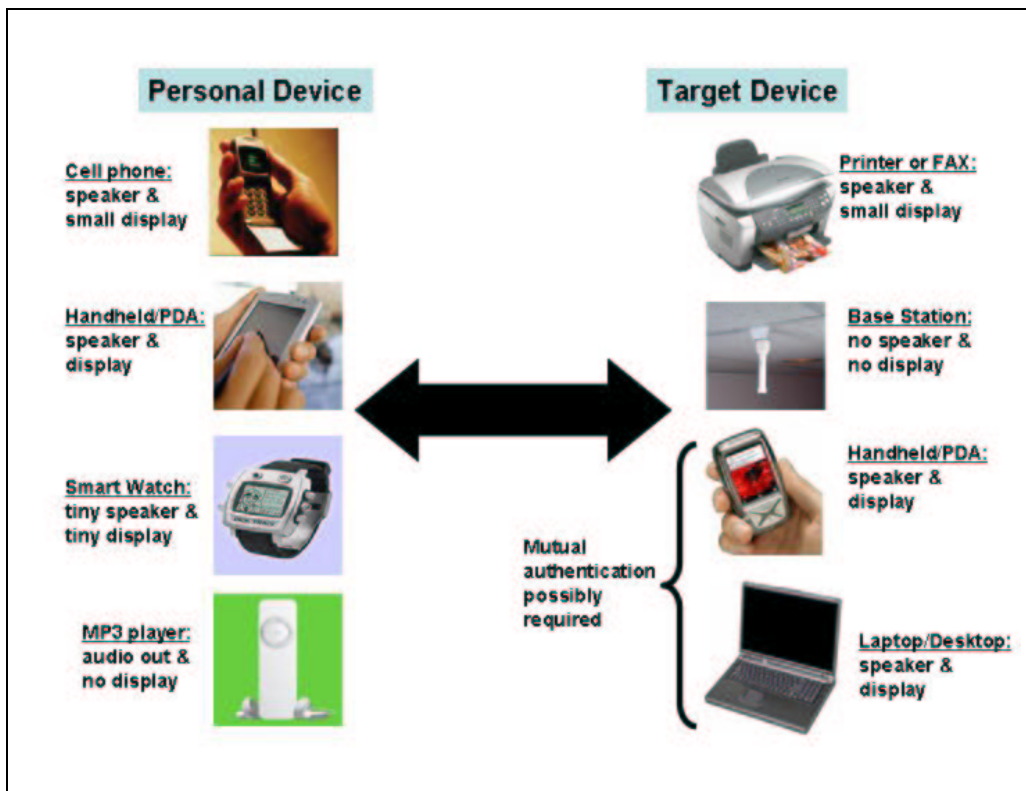


Figure 1. L&C Sample Use Scenarios.

4.2 Classifying L&C Use Cases

Let us consider in more detail the factors that distinguish uses of L&C, including the type of authentication objects, directionality, the number of human users, and the device equipment.

The first factor that distinguishes use scenarios for L&C is the *directionality* of authentication, i.e., whether authentication is one-way (unidirectional from target to personal device) or mutual (bidirectional between the two devices). For example,

in the former, the personal device may need to authenticate the target device's public key and, in the latter, each device may need to authenticate the other's public key authentication object. Since the bidirectional use scenario is a trivial extension of the unidirectional one, in this paper, we focus only on one-way authentication.

A second factor is the number of human users. One classification has a single user with a personal device and the target device does not have an online human administrator. Another classification has two users, each with a personal device. In the former case, if mutual authentication is needed, the burden on the sole human user is essentially doubled.

More than anything else, the equipment available on each device influences the particulars of authentication. Some devices have both a display and a speaker, while others may have only one of these. Thus, some devices, such as low-end base stations, may at first appear unsuitable for L&C, since they lack both a speaker and a display. Nevertheless, they can be accommodated (in a manner similar to the way SiB handles display-less target devices), by affixing a sticker to the target device, which contains the L&C textual encoding of the target device's public authentication object, displayed in print and/or Braille form. In addition, any device, such as a point-of-sale device, with an embedded printer (or a printer itself) is easily supported by L&C, by having the device print an L&C encoding of the authentication object.

We consider four possible use cases for verifiable authentication of some public key:

TYPE 1: hear and compare two audible sequences, one from each device.

TYPE 2: hear an audible sequence from the target device and compare it to text displayed by the personal device.

TYPE 3: hear an audible sequence from the personal device and compare it to text displayed by target device.

TYPE 4: compare text displayed by the personal device to text displayed by target device.

TYPE 1 is the most taxing on the human user of the four; even if the audio sequences are short, comparing them is difficult due to different audio characteristics and the need for the human user to temporarily remember one sequence while waiting to hear the other. However, there is evidence [30] that the two sentences can be compared while being vocalized simultaneously, reducing the user's memory requirement, as well as the total L&C session time.

The TYPE 2 and TYPE 3 use cases are similar, but not the same, since the actual user actions are not identical for these two types. For example, listening to one's own device at very close range is always possible and convenient. Whereas, listening to an unfamiliar target device may require attuning to an unexpected volume, pitch, voice, and noise.

Concerning use case TYPE 4, it clearly does not involve any use of the audio channel and requires non-visually-impaired users. Nevertheless, we include it among our supported cases, since it may be used as an alternative or fall-back method when TYPE 1 to 3 use cases are not plausible. This happens when both devices only have displays, in noisy environments (e.g during a concert) or when silence is required (e.g in a library). If the audio channel is clear, but the devices have displays and no speakers, human users could perform the role of the text-to-speech engine and manually carry a TYPE 2 or TYPE 3 L&C session. Apparently, this use TYPE is not plausible under very low luminance conditions or for visually-impaired users. Hence, L&C is not usable in environments where both the audio and visual channels are unavailable.

Of course, the equipment factor ultimately determines the amount of burden L&C places on the human user. Table 1 shows the types of user requirements corresponding to possible personal-target device combinations. Looking at rows 3 and 6, the choice between TYPE 3 or 1 and TYPE 2 or 1, respectively, can be dictated by the certain properties of the environment. For example, insufficient light, smoke or fog can make TYPE 1 the only viable choice. A visually-impaired user is also likely to choose TYPE 1 over TYPE 3 or 2, unless one of the two devices has a Braille display or sticker. Likewise, the TYPE 4 use case is infeasible for a visually-impaired user unless both devices have Braille displays (or a Braille sticker, in the case of the target device). In row 7, the choice between TYPE 3 or 4 is less clear. One deciding factor might be the comparative quality of the personal device's display and speaker.

4.3 Vocalizable and Readable Representations

In L&C the hash of the target device's public key must be verified by the user(s) of one or both devices. Comparing long (e.g., 160-bit) hashes is a tedious and cumbersome task for the average user. In order to make the process faster and less tedious, a hash must be represented in a more convenient form. In L&C we represent the authentication object as a syntactically-correct text, with the expected use cases being situations (as in TYPE 2 or TYPE 3) where the user reads along with an audio text-to-speech reading of the text.

The text strings generated in L&C are based on the "MadLib" puzzles commonly used by children². That is, we generate a sentence that is a syntactically-correct (but usually non-sensical) from a string of bits.

This string of bits used in L&C will typically be the output from a one-way hash function. Our current implementation allows users to choose between SHA-1, MD2, or MD5 hash algorithms, which is then collapsed into smaller size using a

²In a MadLib puzzle, a funny story is created by having blanks in a text filled in with syntactically-appropriate words chosen by the player or his/her friend.

Row	Use Type	Personal Device		Target Device	
		Display	Speaker	Display	Speaker
1	1	no	yes	no	yes
2	3	no	yes	yes	no
3	3 or 1	no	yes	yes	yes
4	2	yes	no	no	yes
5	4	yes	no	yes	no
6	2 or 1	yes	yes	no	yes
7	3 or 4	yes	yes	yes	no
8	1,2,3 or 4	yes	yes	yes	yes
9	n.a.	no	no	*	*
10	n.a.	*	*	no	no

Table 1. User Requirements for Various Device Combinations. The Use Type column indicates the allowed types of use cases, depending on the device characteristics indicated for the personal and target devices. We use '*' to denote a don't-care condition, we allow the 'Display' condition to include the ability to print or have an affixed sticker attached to device, and we allow the 'Speaker' condition to include any audio-out interface.

technique similar to the “folding” that is done in the S/KEY One-Time Password System [18]. S/KEY partitions an MD5 hash into four 32-bit words, exclusive-ORs the first and third words as well as the second and fourth words. The resulting 64-bit quantity is then used to compute S/KEY phrase-words. L&C collapses the output to a smaller size by first dividing the hash into 10-bit sections. The sections are then paired off sequentially and XOR-ed. The process is repeated for each pair until the desired level of entropy is attained. For example if half of the original entropy is needed, the hash output is collapsed only once. After the hash is collapsed to the desired level, the result is again divided into 10-bit sections. The number of 10-bit sections becomes the final length of the MadLib sentence. (For example, using SHA-1 with 80 bits of entropy would result in 8 MadLib keywords.)

Once the size of the binary input string is determined, an appropriately sized MadLib text is constructed. The text is generated from a template, which consists of a grammatical sentence (or group of sentences) with missing words, each of which being of various types, such as noun, adjective, adverb, verb, boy-name, girl-name, or animal. Each missing word is replaced with a word from a dictionary of appropriate words. The word replacing the MadLib keyword is determined by converting a 10-bit section of the collapsed hash into an integer and using that as the index into the internal dictionary database. For example, the following is a MadLib encoding produced by our prototype, for encoding a 70-bit string (the filled-in words/word-phrases are shown in all caps):

DONALD the FORTUNATE BLUE-JAY FRAUDULENTLY CRUSH-ed over the CREEPY ARCTIC-TERN.

Thus, in order to construct auditorially-robust text sequences, we need to produce a number of word lists of appropriate parts of speech, with the words in each list being as phonetically distant from each other as possible. Using metric for phonetic distance similar to that used by PGPfone [19], but restricted to words of the appropriate type, we can create auditorially-robust word lists for each of the word categories used in our MadLib sequences. In particular, we can do this as follows.

1. Construct a large set C of candidate words of the appropriate type. These should be common English words that can all be used in same place in a MadLib text sequence.
2. Select a random subset W of 2^k words from C , where k is the number of bits we wish to have this type of word represent (e.g., 8 bits for any noun).
3. Repeatedly find the closest pair (p, q) of words in W (using the phonetic distance metric) and replace q with a word from $C - W$ whose distance to any word in W is more than $d(p, q)$, if such a word exists. The resulting set will be a collection of phonetically well-spread words.
4. Order W so that each pair of consecutive words in W are as far from each other as reasonably possible. Doing this optimally is NP-hard [31] but we can use a heuristic algorithm based on pairwise swapping of words in W to come up with a good order for W .
5. Assign integer values to the words in W according to a Gray Code, so that consecutive integers differ in exactly one bit but their respective code words are distant.

This algorithm helps to make the assignment of words to bit sequence auditorially robust—small changes (even to one bit) in

the input should result in noticeably different sounding text strings.

Having described the use types for L&C and the way to produce MadLib text sequences from authentication objects, we describe some example use cases in more detail in the next few sections.

5 Unidirectional Authentication

Providing one-way authentication between devices that have no prior shared context is a challenging problem. Using L&C, the human participates in the procedure and she is able to authenticate the public key of the target device. Using L&C in unidirectional authentication eliminates the need for a trusted party or any pre-shared secret. L&C supports all four types of use scenarios for unidirectional authentication (see Section 3).

In our first example, we consider a target device without a screen or a speaker, e.g., an 802.11 wireless access point. Any wireless device connected to the access point can read out the MadLib sentence generated from the received public key and the user can authenticate the access point by comparing the sentence read out in his machine with the one seen on the sticker attached to the access point.

Another common example is using a printer in a public place. Similar to the above, a printer could have a sticker with the sentence corresponding to its public key. However, a printer is capable of visualizing sentences using its printing functionality. With the help of a button on a printer, it can be asked to print its MadLib sentence and the printer can be authenticated by the identical sentences that are heard from the computer and seen on the paper.

Sound is one of the main elements that enables visually impaired people to interact with the outer world. In this sense, L&C could be easily used to help visually impaired people to authenticate devices or even identities. One basic example could be the authentication scenario of bank ATM's. A visually impaired person could be given a MadLib sentence when she opens a bank account. This MadLib sentence is generated from the card number and its expiration date but using a keyed hash, with a secret key known only by the bank. When the card is inserted in an ATM, the ATM will generate and read the MadLib sentence corresponding to the account number.

6 “Presence” Confirmation

Similar to SiB, L&C can be used to prove the “presence” property of a separate device. A device equipped with a speaker or display can use L&C to confirm, without the assistance of a human, the presence of some other user-attended device within a given range.

To detect the presence of a nearby device, the device generates a nonce value, hashes the nonce value, and converts the hash into the correct MadLib. The device then either displays or automatically speaks out the MadLib, noting the time when the MadLib was displayed or spoken. A user who is in the proximity of the device should be able to read or hear the MadLib and enter the MadLib sentence into their personal device. The personal device can then simply reverse the MadLib encoding process by determining the index of the word in its corresponding dictionary and using the bits of the index as bits of the hash. Once the entire hash has been reverse-computed from the MadLib then the personal device can broadcast the hash of the nonce over the wireless channel. When the original device receives the hash of the nonce over the wireless channel, the device knows that a human user has been present at the display's line-of-sight or within hearing range of the speaker. We assume here that it is hard for a user to strategically place undetected speech or text recognition devices that would relay this information to a remote device.

While the “presence” property can be achieved by devices, it is not robust. It is impossible for the device to determine how many devices are within hearing range or how many are capable of seeing its display. It is also difficult to differentiate between a user in the proximity and a user that has strategically placed speech or text recognition devices that would relay the audio or visual signal to a remote device. However we assume, that it is hard and not beneficial enough for a user to do so undetected. The device is only capable of verifying that the hash received over the wireless channel is the same as the one generated. If desired, it is possible to determine the time difference between when the nonce hash is received over the wireless channel and when it was first displayed or spoken.

Even though the “presence” property is not robust, it is not without its own uses. A device can use the “presence” property to limit the authority to control the device to users located within view or hearing range of the same device. This property can also be used as a means of triggering the power saving feature of a device, when no user-attended device is within visual or audio range. Before a device switches to a lower power consumption state, the device displays or speaks the MadLib. If no response is received over the wireless channel after a certain time period, the device assumes that it is not needed by a user-attended device in the proximity and that it can safely switch state.

7 Implementation and Performance

In this section we describe the prototype implementation of L&C and experimentally assess and analyze its performance, taking into account various cost factors.

7.1 Implementation Overview

Since L&C is intended for a variety of mobile computing platforms, portability is a key requirement. We built L&C using the highly portable *Ewe* Programming System [32] which facilitates the development of Java applications.³ Our implementation runs on any Pocket PC (iPAQ in our experiments) and any Windows PC.

As part of its initialization, L&C allows selection between bi-directional and uni-directional public key authentication. In the uni-directional case, L&C runs on Alice (a user-attended personal device) to verify the public key of Bob, an un-attended target device, as discussed in Section 3. In the bi-directional case, both devices are user-attended, either by a single user or two different users. Also at bootstrap time, L&C allows the choice of 802.11 or Infrared communication channel. One of the participants/devices (say, Bob, the target device) initializes L&C and waits for a TCP connection on a well-known port over one or both communication channels. Alice's user physically approaches Bob (particularly relevant for Infrared communication), initializes its L&C application and connects to Bob. Next, Bob sends its public key to Alice via the selected communication channel. In case of bi-directional authentication, Alice reciprocates. On Bob, L&C converts the hash of the local public key into a MadLib sentence (as discussed above) and displays sentence. On Alice, L&C converts the hash of the received public key and displays the corresponding MadLib sentence. In the bidirectional case, in addition to the above, Bob converts the hash of Alice's public key and Alice generates the MadLib of her key. The user(s) have the option to vocalize the sentences on both devices.

Recall that there are four types of user requirements or settings. In Type 1 setting, the user hears the same MadLib vocalized by both devices. In Type 2 and Type 3 settings, the user reads the MadLib from one device and compares it to the same MadLib vocalized by the other devices. Finally, Type 4 involves the user reading and comparing two Madlibs, one displayed by each device.

Although the main purpose of L&C is to authenticate one (or both) public keys, this process usually serves as a prelude to a key exchange protocol, i.e., the generation of a session key to be used for subsequent secure communication between the two devices in question. For this reason, L&C includes two flavors of public-key-based key exchange protocols: RSA- and DH-based. These are basically standard textbook protocols (e.g., the Station-to-Station protocol) and we do not elaborate on them further.

7.2 Implementation Details

Owing to modular software design, L&C can utilize a variety of Text-to-Speech (TTS) engines. However, most C/C++ speech engines are platform-dependent, while those written for mobile devices are mostly proprietary. Furthermore, Java-based TTS engines are available for specific JVM-s that are unsuitable for resource-constrained devices, such as smartphones and iPAQ-s. Specifically, Sun offers the JSAPI and FreeTTS Java TTS engine implementations. However, these run only on Java 1.4. Therefore, we employed existing TTS applications that could be used by L&C—*Digit for PC* and *Pocket PC* by Digalo⁴, which is a simple lightweight clipboard reader that uses the *Elan Speech Engine*.⁵ Our application copies the text to-be-vocalized onto the system clipboard and Digit speaks it out automatically or when the user presses a button on the application window. Digit is initialized and terminated from within the Ewe program.

Ewe does not provide a complete API for low-level cryptographic primitives. Thus, in order to implement DH- and RSA-based key exchange protocols, we added a lightweight cryptographic API to Ewe's Java libraries. For this purpose, we ported the *Bouncy Castle* crypto package [33] for JDK 1.3. For hashing we used Ewe's built-in SHA-1.

The FreeTTS and Bouncy Castle crypto package are written solely in Java and do not link to native platform-specific libraries, facilitating L&C's platform independence. So far, we tested L&C on Pocket PC and Windows PC. L&C can also be used with the rest of the Ewe-supported platforms platforms by changing the TTS engine. We are currently porting Sun's FreeTTS and JSAPI into Ewe.

NOTE: L&C source code, installation instructions as well as pictures and a video demonstrating L&C can be found at www.ics.uci.edu/ccsp/lac.

³Ewe is currently available for the following platforms: Pocket PC (Windows CE), MS SmartPhone, Casio BE-300, HandHeldPC Pro, Sharp Zaurus, Linux PC, Windows PC and any Java 1.2 VM.

⁴See: www.digalo.com.

⁵See: www.elantts.com.

7.3 L&C Performance

In this section we evaluate L&C performance using a commodity laptop PC as a target device and a low-end iPAQ as a personal device. The laptop PC is equipped with: Intel Pentium M Centrino 1.7GHz, 400MHz FSB, 2MB Cache and 512 MB RAM running Windows XP. The iPAQ is a Compaq 3650 equipped with: an Intel Strong ARM SA-1110 32-bit RISC processor operating at 206MHz, with 31.25 MB RAM, 16 MB ROM running Windows CE version 3.0.9348 (Pocket PC 2002). For the 802.11g channel we configured a wireless subnet consisting of: one wireless router, two iPAQs and a PC. The channel's nominal bandwidth for all devices is 54 Mbps. The Infrared ports of all devices operate at 115 Kbps.

As mentioned above, once L&C key authentication completes and after the user or users have verified each other's public keys, the two devices proceed with establishing a shared secret. However, the protocol for generating a shared secret does not require any human involvement. Therefore, it is omitted from the following performance analysis.

We analyze L&C performance for human-verifiable authentication of either Diffie-Hellman (DH) or RSA public keys. The system-wide known DH parameters (p , g and q) and the RSA public exponent e (see Appendix) are neither sent nor verified by L&C. Furthermore, the DH key pair and the RSA public key are generated off-line and do not contribute to protocol completion time, regardless of whether they are ephemeral or long-term. L&C can generate a new DH key pair for the same DH parameters in 3540.2 and 272.1 ms on the iPAQ and PC, respectively. The corresponding times for the generation of RSA key pairs is significantly larger, since prime number generation is involved. Note that we use 1024-bit moduli for both RSA and DH (see Appendix).

Table 2, lists processing times only for the operations that involve (or lead to) human-verifiable authentication of public keys. Tables 3(A), 3(B), and 3(C) show timings for different types of L&C unidirectional authentication sessions. The corresponding bidirectional sessions can be analyzed in a very similar manner, therefore we do not include it in our analysis. Operations are listed in the order they take place.

Measurements for operations 1 through 5 are obtained as the average over 20 L&C sessions operated by human users. The times for other operations are obtained over 300 bulk repetitions of L&C sessions that do not include the above four operations. We use Ewe's timing function, which offers (only) 10 ms precision. L&C initialization times (row 1) are obtained after RAM has been reset, so that they include the time to load all the application (Ewe VM, L&C class files and Digit) into memory. The total time does not reflect any further delays introduced by a human user, hence, it represents a rough approximation.

No	Operation	iPAQ	Laptop PC
1	Ewe VM and GUI initialization	2430	120
2	Digit initialization	18310	1092
3	L&C setup by user	1502	910
4	TCP Connection est., 802.11g	3.2	0.4
5	TCP Connection est., IR	3.4	-
6	pub. key transmission, 802.11g	5.6	0.1
7	pub. key transmission, IR	6.1	-
8	pub. key MadLib generation	365.6	17.1
9	pub. key MadLib vocalization	4791	4637

Table 2. Average processing times (in ms) of L&C operations.

Ewe VM and GUI initializations take place while Digit is initialized. For L&C running on a PC, these initializations complete almost simultaneously, allowing the user to proceed with L&C setup while Digit bootstraps. On an iPAQ, Digit initialization preempts the processor, not allowing the user to use the GUI to setup the session. Therefore, at this phase of the L&C session, Digit initialization is the only operation that needs to be considered.

In reality, the time a user spends on L&C setup is comprised of: (1) entering the target device's network address (IP address or "infra-red"), pressing "Enter" and "Connect" buttons and, lastly, aligning the devices if the Infrared channel is used. It also includes the time for the application to reach the *accept()* or *connect()* calls. For the experiments, the Infrared ports were pre-aligned and default network addresses were used. Hence, the user needs only to press two buttons to initiate the connection.

Connection establishment time (operation 4 or 5 in Table 2) is the time required by the TCP socket *connect()* system call to connect to the accepting process running on the iPAQ or the PC. (We measured L&C sessions over the IR channel only

between iPAQs.) Times for operations 10 and 11 vary with the length of the Madlib sentence.

Operations 1, 2, 3, 6, 7 and 8 in Table 2, take place concurrently on both devices. Therefore, only the lengthier of the two counts towards the total time. We do not measure the time for reading the Madlib sentence from the device’s display since it is user-dependent. In all experiments, we use syntactically-correct MadLib sentences consisting of 10 words, of which 7 are S/KEY-generated. The sentence format is the same with the one presented in Section 4.3.

Initialization	22,245
----------------	--------

(A) Type 1

Table 2, Row 6 Col. 3	5.6
Table 2, Row 8 Col. 3	365.6
Table2, Row 9 Col. 4	4,791
Table 2, Row 9 Col. 3	4,637
TOTAL TIME	32,044

(B) Type 2

Table 2, Row 6 Col. 3	5.6
Table 2, Row 8 Col. 3	365.6
Table 2, Row 9 Col. 4	4,637
TOTAL TIME	27,253

(C) Type 4

Table 2, Row 6 Col. 3	5.6
Table 2, Row 8 Col. 3	365.6
TOTAL TIME	22,616

Table 3. Timings (in ms) for TYPE 1, 2 and 4 L&C Sessions.

Table 3 shows the timings for TYPE 1, 2 and 4 L&C sessions. The row labeled “**Initialization**” at the top of the table reflects the sum of rows 1, 2, 3 and 4 (iPAQ column) of Table 2. This significant delay – almost 22 seconds – is induced by all the non-cryptographic software initialization on the iPAQ. Also, this delay is independent of the L&C use case.

Table 3(A) examines the absolute worst case scenario under unidirectional authentication – that involving a user performing one-way, voice-only (TYPE 1) authentication of the target device’s public key. The user verifies the target device’s (PC) public key by hearing the corresponding MadLib spoken by the PC and its iPAQ device, in either order. (Our preliminary user studies indicate that it is preferable for MadLibs to be vocalized simultaneously; otherwise, users have difficulty understanding the sentences.) This means that our scenario includes two MadLib generations: one per device which can take place concurrently. In addition, it includes two MadLib vocalizations, which must take place sequentially. The row labeled “TOTAL TIME” reflects the sum of the times for all individual operations. As the results show, the TYPE 1 unidirectional session between the iPAQ and the laptop PC can complete in approximately 32 seconds.

Table 3(B) examines the most commonly anticipated use scenario, which corresponds to TYPE 2 of Section 3. It involves uni-directional audio-based authentication of the target device’s public key. The user-attended iPAQ receives the public key of the target device (PC), hashes it and generates the corresponding MadLib. Then, the user reads the MadLib sentence from the iPAQ’s display and compares it to the vocalization by the PC. As can be seen in the table, the time required by this type of L&C session totals approximately 27 seconds.

Table 3(C) shows timings for a uni-directional display-only (TYPE 4) L&C session, assuming 802.11g channel. The actual time would include that needed by a user to read and compare the two displayed (one on each device) MadLib sentences. We did not measure TYPE 3 as it is very similar to TYPE 2.

7.4 Performance Analysis

Table 2 clearly illustrates that the overall cost is dominated by: (1) Ewe VM and GUI initialization and (2) MadLib vocalization. The first can be omitted if multiple L&C sessions take place after a single initialization or if L&C is pre-initialized. Naturally, the time to speak out a MadLib sentence is proportional to the word length of the sentence as well as to the number of syllables in each S/KEY-generated word. Therefore, it can vary for the same word-length sentences.

MadLib generation is approximately 20 times more expensive on Pocket PC. It is also evident that the time for a typical user to set up L&C is greater on an iPAQ than on a PC. This is due to the rather user-un-friendly GUI and slower rendering of the GUI components on the Pocket PC.

Communication costs represent only a tiny fraction of the total cost. Transmission of the public keys is an order of magnitude cheaper on a PC. Since packet traversal of the protocol stack – and not the actual transmission over the physical medium – is the dominating factor of the communication cost, the difference in communication costs over the 802.11g and Infrared channels is truly insignificant.

Processing and memory limitations on the iPAQ result in significantly longer delays than on the PC. However, we stress that we used old and **low-end** iPAQs in our experiments. Using current state-of-the-art iPAQs or other similar devices would greatly reduce delay. Experiments suggest that the most plausible way to reduce protocol overhead is to shorten MadLib sentences and speed up TTS engine initialization time. From the above analysis, we conclude that L&C is a viable solution on platforms with moderate computation and communication capabilities.

7.5 Performance Improvements

A US patent addressing a similar problem to ours [24], proposed augmenting the DH key agreement with additional steps. Upon completion of Diffie-Hellman exchange, both devices one-way hash the agreed-upon session key. Then they truncate the hash bit sequence to the desired length by taking its leftmost bits. The device(s) can display this bit sequence (thereafter referred to as check value) allowing the user(s) to compare them and verify that both devices have a common session key.

The pitfall of hash truncation lies with the number of digits that need to be examined by the user. Reading or typing a large number of digits using the small interfaces of a personal device without error can be cumbersome. Drastically truncating the hash-code, e.g. to the first 16 or 32 bits results to a serious security weakness.⁶ Based on the assumed adversarial capabilities (see Appendix and Section 8.1), $t = 50$ and $t = 70$ provide the necessary security when ephemeral and one-year-term public keys are used, respectively.

When ephemeral DH public keys are used, examining 13 hexadecimal digits is an error prone task for a human. L&C renders such comparison user-friendlier, because the user compares five-S/KEY-word Madlibs. The time to vocalize a five-S/KEY-word sentence is on average 2905 ms on a PC and 3096 ms on the iPAQ. Shorter MadLib sentences yield reduced vocalization time and easier comparison for the user. Indicatively, the TYPE 2 scenario described in Table 3(B) would require a total of approximately 25.5 seconds to complete.

In [25] the proposed solution prevents attackers from finding second pre-images for long-term public keys, without requiring the users to examine too long hexadecimal sequences. The schemes of interest are called MANA-I and II and they employ keyed check-functions that use short (10-20 bit) keys and produce short check-values. These check-functions are essentially MAC (Message Authentication Code) functions. We denote this check function as $MAC_{k,t}(V)$, where k is a random key, t is the length of the check value and V is the verification object (e.g a public key). The MANA-I-inspired L&C would operate as follows in the case of unidirectional long-term public key authentication:

1. Bob sends Alice his public key (e.g g^b) using the unprotected wireless channel.
2. Bob generates the random 20-bit key k to use with the check-function. Bob also generates the check-value $MAC_{k,50}(g^b)$. Bob generates the MadLib sentence for the check-value (five words) and a 4 digit hexadecimal sequence for the random key. Bob displays or vocalizes the MadLib and the 4-digit sequence.
3. The user enters the presented by Bob random key to the device Alice.
4. Alice uses the random key to recompute the check-value on the received g^ab , and displays or vocalizes the check value MadLib.
5. The user completes the process by comparing the values displayed by the two devices. Only if the check-value Madlibs are the same, the exchange is accepted by the user.

Using MANA-I, the length of the MadLib sentence, when long term public keys are used, is reduced from the recommended seven S/KEY words to five. However, the user is required to enter a four digit hexadecimal code in one of the two devices, using a keyboard or a touchscreen, certainly resulting to more time-consuming verification.

To avoid entering the key, we could use MANA-II, in which the random key is transmitted over the wire (becomes available to the attacker) and is displayed for comparison on both devices. For long term public key verification, it would still require the user to compare a five S/KEY word sentence plus two S/KEY words for the key, yielding no actual gains. However,

⁶The attacker can replace Alice's transmitted g^a value with $g^{a'}$, wait for Bob to reply with g^b and intercept his transmission. Next, the attacker can perform exhaustive search to find $g^{b'}$ such that $h_t(g^{a'b}) = h_t(g^{a'b'})$, and relay $g^{b'}$ to Alice as Bob's public key. On average, the attacker needs to perform 2^{t-1} modular exponentiation and hashing operations before a suitable b' value is found.

MANA-II can be used to authenticate the agreed upon shared secret, when long-term public keys are used, without exposing substantial information about it. The use of the random key further masks information about the secret.

In [28], Cagalz et al. proposed the *DH-SC* protocol for human-verifiable authentication. This protocol uses a commitment scheme, which transforms a value m into a commitment/opening pair (c, d) . In this pair, c reveals no information about m (e.g c is public key encryption of m) but c and d (e.g d is the encryption key) reveal m . In an ideal commitment scheme it is infeasible to find d' such that (c, d') open to $m' \neq m$. Their protocol requires the communicating parties to compare a string derived from the XOR of two per-session random bit-sequences contributed by the two parties. Hence, unlike MANA-II, the users do not have to compare a value (the random MAC key), which does not contribute to the uncertainty of the attacker. Hence, *DH-SC* effectively reduces the length of the compared values for a given level of security. For example, two-S/KEY-word MadLibs generated from the random bit-sequences provide security almost equivalent to the one of five-S/KEY-word MadLibs derived from the hash of the ephemeral DH public keys (see Section 8.1).

The *DH-SC* protocol proceeds as follows: Both Alice and Bob (A and B) generate their public keys g^a and g^b , respectively. Then A and B each generate a t -bit random string N_A and N_B . They use them to calculate commitment/opening pairs for the concatenations $0\|g^a\|N_A$ and $1\|g^b\|N_B$ (0 and 1 are fixed values used to prevent reflection attack). In the first message, A sends to B the commitment c_A and B responds with his commitment c_B . In turn, A sends her commitment key d_A with which B opens c_A and obtains $g^{a'}$ and $N_{A'}$. B checks the correctness of the commitment pair c_A, d_A and verifies that 0 appears at the beginning of the message. If the verification is successful, B sends d_B , with which A opens c_B and obtains $g^{b'}$ and $N_{B'}$. A checks the correctness of the commitment and if it is valid, both parties proceed with generating the verification strings for $i_A = N_A \oplus N_{B'}$ and $i_B = N_B \oplus N_{A'}$, respectively. Now the users of A and B can simply compare the verification strings (MadLib sentences for L&C) and accept the exchanged public keys only if they match.

8 Security Analysis

Assuming that all underlying cryptographic primitives, cryptographic hash functions and public key algorithms (RSA and Diffie-Hellman key agreement) are secure, the security of L&C depends on the adversary's inability to: (1) interfere with the audio or visual channel and (2) compromise the target device. In this section, we discuss the security of the cryptographic primitives and compare the security of the alternative channels. We then discuss some concrete attack scenarios.

8.1 Cryptographic Primitives

To establish a secure bidirectional channel between communicating devices (as described above), L&C uses either the ephemeral DH key agreement or RSA-based mutual key agreement (see Appendix). The security of both methods is based on well-known and believed-to-be-hard computational problems. Both types of key agreement can be protected against man-in-the-middle (MITM) attacks by verifying the exchanged public keys. Security against MITM attacks can also be achieved by using L&C to verify the hash of the agreed-upon shared secret and a random value known by both parties. However, public key verification is preferable because the users can detect the attack early in the process and abandon it without performing further expensive computations. Also it is common for users to exchanged public keys in order to establish secure communication at a later time, in which case they do not need to generate the shared secret immediately. Furthermore, it is wise not to expose any information about the shared secret by using L&C to verify it.⁷

In the DH case, the two devices exchange their respective public keys. By verifying that the received DH public key originates with the intended sender, the user(s) can be certain that the DH session key is the same for the two parties.

In one variation of the RSA-based scheme, the two devices exchange their public keys, which they use to encrypt session key contributions. If the user(s) verifies, using L&C the other party's public key, he can be certain that the encrypted session key contributions cannot be decrypted by an adversary and the final session key is secure. Even if the adversary interferes with the transmission of the encrypted key contributions, the result would be failed key agreement which can be easily detected by both parties.

The devices can optionally use their RSA keys to sign a well-formatted message containing the encrypted partial secrets. In this way they ensure that the encrypted partial secrets originate with the intended sender and avoid the delay involved in detecting shared secret inconsistency if an active attack were to occur. Alternatively, they can use displayed or vocalized Madlibs to verify the established shared secret. Although shared secret verification is sufficient to detect MITM attack at any stage of the protocol, they can combine it with public key verification to possibly detect an attack earlier.

In L&C the hash function of the transmitted public key is converted to a MadLib sentence. During the conversion, the length of the MadLib sentence can be set to the desired entropy level. If the adversary can find a second pre-image of the

⁷Unless we employ a secure MAC function as is the case with the MANA-II-inspired protocol presented in the previous section.

hash value that encoded in the MadLib sentence, then it can attack the communication channel and replace the transmitted key with the collided pre-image, causing the two devices to falsely verify the key agreement. However, second pre-image resistance is a required property for cryptographic hash functions and we assume that it holds.

Although the length of the MadLib sentence can be set to achieve the desired level of security, it is commensurate with the burden of comparing longer sentences. In many cases, the user is not expected to memorize the sentence, but a longer sentence still results in a longer verification time. The vulnerability of using short sentences becomes real when long-term DH public keys are used. In this case, the adversary may have enough time to find a second pre-image via a brute force attack.

According to our experimental results, the DH public key exchange in L&C requires only a few milliseconds (around 6 ms) even in resource-constrained PDAs. The attacker would need to perform on average 2^{t-1} modular exponentiation and hashing operations to determine a' such that $h_t(g^{a'}) = h_t(g^a)$, where $h_t()$ denotes the collapsed to t bits one-way hash. For $t = 30$ the attack requires roughly 530 million trials. Under our adversarial model (see Appendix) and assuming that the key exchange “times out” after at most 100 ms, the attacker is able to perform 100 million trials. Therefore, $t = 50$ (five-S/Key-word MadLib) should suffice to ensure that, in the case of ephemeral DH, the attacker can succeed only with probability roughly equal to 2^{-22} . Once a session key is established, finding the collision on the MadLib sentences does not give the adversary any advantage. The public key generation for RSA is significantly costlier rendering the above attack infeasible even for smaller values of t .

The attacker may resort to off-line pre-computation of public keys that hash to all possible t -bit hash values, opting to derive second pre-images by an efficient table lookup at key exchange time. He can compute the MadLibs for 2^t public keys and store these values. He can repeat this process with new public keys until he derives second pre-images for more t -bit hash values increasing his probability of success. Taking into account the indexing information and that table entries consist of the t -bit hash values and at least two public keys that hash to the same value, the storage required is more than 2^{t+8} bytes. Based on our assumptions (see Appendix), a powerful attacker can perform roughly 2^{46} trials per year, thus for $t = 50$ he can succeed with significant probability after a year of computations. The actual limitation is the storage requirements which would be as high as eighteen petabyte. In order for the attacker to succeed with probability higher than 2^{-20} , he would need to store more than 2^{38} bytes. Still, the table lookup operation would be unlikely to complete in a few milliseconds, even if we assume very efficient indexing and lookup algorithms and very fast disk and memory access times.

Given the assumed adversarial capabilities, if **one-year-term** DH public keys are used, $t = 67 - 70$ (seven S/Key generated words) is sufficient to prevent an attacker from finding second pre-images with significant probability.

If L&C uses the *DH-SC* protocol, generating MadLib sentences from the verification bit-strings i_B, i_A , the attacker can succeed only by guessing correctly the random values N_B, N_A . He can use N_B, N_A to create commitments for his public keys, which can in turn be used to mount an active attack. Since he has only one opportunity to send fake commitments, the probability of success, given $|N_B| = |N_A| = t$, is 2^{-t} . Hence, $t = 20$ (two-S/KEY-word MadLib) provides security roughly equivalent to the one achieved using five-S/KEY-word MadLib derived from the hash of the ephemeral DH public keys. The security of *DH-SC* is discussed thoroughly in [28].

8.2 Alternative Channels

Authentication of communicating parties without the assistance of a trusted party requires a communication channel that is secure against active attacks. In this section, we discuss several alternative channels and their security compared to that of L&C.

Wired channels between two devices, such as cross-over Ethernet, direct USB or serial connection, can be considered secure against active attacks since both ends of the cable are directly connected to the intended devices and the entire cable is visible to the human user. However, lugging around different types of wires in order to communicate with different types of devices is quite cumbersome.

Wireless alternatives, such as 802.11a/b/g, Infrared and Bluetooth, are not secure against MITM attacks. A human user is essentially incapable of identifying which two devices are actually communicating over a wireless channel. The only way to be sure that the intended devices are communicating without interference is through status indicators on devices, which is prone to errors and attacks. Among wireless alternatives, Infrared is considered most secure since it requires a direct line-of-sight between devices (which can be visually inspected). However, it is still possible for a well-placed adversary in a device-crowded setting to interfere with communication, since the human user cannot determine the exact direction of an infrared signal.

As discussed in Section 1, the visual channel – exemplified by SiB is a very viable alternative. However, many devices lack cameras as well as large-enough displays and powerful-enough CPUs to perform necessary image processing (e.g., to

extract the barcode). Furthermore, environmental factors inhibiting the use of cameras as well as sight-impaired users limit the applicability of the visual channel. Other hash visualization techniques [34, 35, 16, 13] are also viable alternatives for verifying public and/or session keys. They can detect MITM attacks as long as the comparison is done very carefully and the environment is well-illuminated.

8.3 Attacks Against L&C

Using L&C, the man in the middle adversary is easily detected. When the personal device’s user asks for the target device to read its MadLib sentence loud, if she is far enough, an adversary device close to the target device could impersonate the target device. A human user however, can trivially determine whether it is the target device talking by checking the screen or sensing the origin of the sound. Furthermore, the user(s) can detect the active attack if the adversary generates sound without previously asking the device(s) to do so.

An adversary may attempt to attack L&C by using malicious software surreptitiously installed on the target device. Malicious software might tamper unnoticed with the input and output of the desired application. This problem is sometimes referred to as *secure windowing* or *establishment of a trusted path to system window* [36]. Such attacks can be prevented by using a good system window policy, e.g., displaying application windows in predetermined locations or, as proposed in [37], using predetermined window backgrounds known only to the user. A malicious application may also try to process the MadLib sentence from screen pixels and redisplay it, which would appear legitimate to the user. This can be prevented by limiting the screenshot and copy operations to the application that is on top at the time. EROS, a trusted window system, proposed by Shapiro, et al. [36] provides further details on how to achieve this type of protection.

An adversary may also emit interference in the form of ambient noise in order to alter L&C audio. However, this attack can be easily detected by a human user since two different sound sources would be involved. Moreover, the adversary would need near-perfect synchronization, a-priori knowledge of the MadLib sentence and plenty of computational power to be able to produce correct sounds that will make the sentence sound like a MadLib sentence corresponding to some public key picked by the adversary.

Excessive ambient noise could also prevent the user(s) from hearing each other’s devices, but the source of the noise would be obvious. Even if it cannot be avoided, ambient noise would at best result in denial of service. In such cases, user(s) can fall back to comparing MadLib sentences displayed or printed on the screens (or stickers) of the respective devices.

9 Conclusions

This paper presented the design of the **Loud-and-Clear (L&C)** system for human-assisted device authentication. L&C places relatively little burden on the human user, since it is based on the audio channel and uses a text-to-speech engine to read an auditorially-robust, syntactically-correct sequence derived from an authentication string. We also discussed some anticipated common use cases and provided experimental performance data for our L&C prototype.

References

- [1] Loren M. Kohnfelder., “Towards a practical public-key cryptosystem.” 1978, b.sc thesis, MIT Department of Electrical Engineering.
- [2] P. Zimmermann, *The Official PGP User’s Guide*. MIT Press, 1995.
- [3] S.P. Miller, C. Neuman, J.I. Schiller, and J.H. Saltzer., “Kerberos authentication and authorization system.” in *Project Athena Technical Plan*, 1987, p. pg section E.2.1.
- [4] Jennifer G. Steiner, Clifford Neuman, and Jeffrey I. Schiller., “Kerberos: An authentication service for open network systems.” March 1998, Manuscript.
- [5] W. Diffie and M. E. Hellman., “New directions in cryptography.” in *IEEE Transactions on Information Theory*, November 1976, pp. IT-22(6):644–654.
- [6] Alfred J. Menezes, Paul C. van Oorschot, and Scott A. Vanstone., *Handbook of Applied Cryptography*, ser. CRC Press Series on Discrete Mathematics and its Applications. CRC Press, 1997.
- [7] D. Kugler, “Man in the middle attacks on bluetooth,” In *Financial Cryptography ’03*, Long Beach, 2003. Lecture Notes in Computer Science, Springer-Verlag.
- [8] M. Jakobsson and S. Wetzel, “Security weaknesses in bluetooth,” *Lecture Notes in Computer Science*, vol. 2020, pp. 176+, 2001.

- [9] Steven Bellovin and Michael Merrit., “Augmented encrypted key exchange: a password-based protocol secure against dictionary attacks and password file compromise.” in *First ACM Conference on Computer and Communications Security CCS-I*, 1993, pp. 244–250.
- [10] V. Boyko, P. MacKenzie, and S. Patel., “Provably secure password authentication and key exchange using diffie-hellman.” *Advances in Cryptology–EUROCRYPT*, vol. 1807 of Lecture Notes in Computer Science, pp. 156–171, 2000.
- [11] P. MacKenzie, S. Patel, and R. Swaminathan., “Password authenticated key exchange based on RSA.” in *Advances in Cryptology – ASIACRYPT*, 2000, pp. 599–613.
- [12] T. Wu., “The secure remote password protocol.” in *Network and Distributed System Security Symposium*, February 1999.
- [13] Ian Goldberg., “Visual key fingerprint code.” 1996, Available at <http://www.cs.berkeley.edu/iang/visprint.c>.
- [14] Ralph Levien, “PGP snowflake,” 1996, source code available at: <http://packages.debian.org/testing/graphics/snowflake.html>.
- [15] Adrian Perrig and Dawn Song., “Hash visualization: A new technique to improve real-world security.” in *Proceedings of the 1999 International Workshop on Cryptographic Techniques and E-Commerce (CrypTEC '99)*, July 1999, pp. 131–138.
- [16] C. Ellison and S. Dohrmann, “Public-key support for group collaboration,” *ACM Trans. Inf. Syst. Secur.*, vol. 6, no. 4, pp. 547–565, 2003.
- [17] Jonathan M. McCune, Adrian Perrig, Michael K. Reiter, “Seeing-Is-Believing: Using Camera Phones for Human-Verifiable Authentication,” in *2005 IEEE Symposium on Security and Privacy*, 2005, pp. pp. 110–124.
- [18] N. Haller, “Available at RFC1760: The S/KEY One-Time Password System,” February 1995.
- [19] P. Juola and P. Zimmermann, “Whole-word phonetic distances and the pgpfone alphabet,” in *Proceedings of the Fourth International Conference on Spoken Language Processing (ICSLP)*, vol. 1, 1996, <http://www.asel.udel.edu/icslp/cdrom/vol1/005/a005.pdf>.
- [20] M. J. Atallah, C. J. McDonough, V. Raskin, and S. Nirenburg, “Natural language processing for information assurance and security: An overview and implementations,” in *Proceedings of New Security Paradigm Workshop*, 2000, pp. 51–65.
- [21] Frank Stajano and Ross Anderson., “The resurrecting duckling: Security issues for ad-hoc wireless networks.” in *Security Protocols, 7th International Workshop*, 1999.
- [22] Dirk Balfanz, D.K. Smetters, Paul Stewart, and H. Chi Wong., “Talking to strangers: Authentication in ad-hoc wireless networks.” in *Symposium on Network and Distributed Systems Security (NDSS '02)*, February 2002.
- [23] S. Capkun, J. Hubaux, and L. Buttyan., “Mobility helps security in ad hoc networks.” in *Proceedings of the ACM Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc 2003)*, June 2003.
- [24] D. P. Maher, “Secure communication method and apparatus,” U.S. Patent Number 5,450,493 September 1995. Filed on 29th December 1993.
- [25] C. M. C. Gehrman and K. Nyberg, “Manual authentication for wireless devices,” *RSA Cryptobytes*, Vol. 7, No. 1, pp. 2937, 2004.
- [26] C. Gehrman and K. Nyberg, “Security in personal area networks,” In C. J. Mitchell, editor, *Security for Mobility*, pages 191230. IEE, London, 2004.
- [27] J.-O. Larsson, “Higher layer key exchange techniques for bluetooth security,” Open Group Conference, Amsterdam, October 2001.
- [28] M. Cagalj, S. Capkun, and J. Hubaux, “Key agreement in peer-to-peer wireless networks,” *Proceedings of the IEEE (Special Issue on Cryptography and Security)*, 2006.
- [29] S. Vaudenay, “Secure communications over insecure channels based on short authenticated strings.” in *CRYPTO*, 2005, pp. 309–326.
- [30] D. L. Greg Kochanski and C. Shih, “A reverse turing test using speech,” in *Proceedings of the Seventh International Conference on Spoken Language*, September 2002, denver, Colorado. [Online]. Available: <http://kochanski.org/gpk/papers/2002/rtt.pdf>
- [31] M. R. Garey and D. S. Johnson, *Computers and Intractability; A Guide to the Theory of NP-Completeness*. New York, NY, USA: W. H. Freeman & Co., 1990.
- [32] M. Brereton, “Ewe java vm for pocketpc,” available at <http://www.ewesoft.com/>.
- [33] “Bouncy Castle Crypto APIs,” Available at: <http://www.bouncycastle.org/index.html>.
- [34] P. D. MacKenzie, S. Patel, and R. Swaminathan, “Password-authenticated key exchange based on rsa,” in *Proceedings of the 6th International Conference on the Theory and Application of Cryptology and Information Security*. London,

UK: Springer-Verlag, 2000, pp. 599–613.

- [35] A. Perrig and D. Song, “Hash visualization: A new technique to improve real-world security,” in *Proceedings of the 1999 International Workshop on Cryptographic Techniques and E-Commerce (CrypTEC '99)*, July 1999, pp. 131–138.
- [36] J. S. Shapiro, J. Vanderburgh, E. Northup, and D. Chizmadia, “Design of the eros trusted window system.” in *USENIX Security Symposium*, 2004, pp. 165–178.
- [37] R. Dhamija and J. Tygar, “The battle against phishing: Dynamic security skins,” in *SOUPS '05: Proceedings of the 2005 symposium on Usable privacy and security*. ACM Press, 2005, pp. 77–88.
- [38] W. Mao, *Modern Cryptography: Theory and Practice*. Prentice Hall Professional Technical Reference, 2003.
- [39] A. K. Lenstra and E. R. Verheul, “Selecting cryptographic key sizes,” *Journal of Cryptology: the journal of the International Association for Cryptologic Research*, vol. 14, no. 4, pp. 255–293, 2001. [Online]. Available: citeseer.ist.psu.edu/article/lenstra01selecting.html

A: Threat Model and Cryptographic Assumptions

Below we discuss our assumptions about the threat model and cryptographic operations used in L&C.

We adopt the Dolev-Yao threat model [38], under which the attacker controls the communication channel; he can acquire any message transmitted over the radio channel and can send messages to any user sharing the channel. We assume however, that the adversary is computationally bounded. Specifically the attacker can perform at most one million modulo a 1024-bit number exponentiations and hashing operations per second. We further assume that two authenticated parties involved in the communication trust each other (e.g an authenticated party will not disclose secret information to a third party).

All Diffie-Hellman key agreement operations are performed over a multiplicative group G , which is a subset of Z_p^* , where p is a large prime. g is the generator of G with order q , where q is a large prime s.t. $q \mid p - 1$. A DH key pair is generated as $(g^a \bmod p, a)$, where $a \in Z_q^*$ is the secret key corresponding to the DH public key g^a . The Decisional Diffie-Hellman problem is assumed to be hard in G .

The RSA-based key agreement operations are performed over Z_n^* , where n is a composite such that $n = p \cdot q$ where p, q are large primes. The RSA public key is n, e and the secret key is d such that $e \cdot d = 1 \pmod{n}$. The factorization of n is known only to the holder of the secret d and is assumed to be hard.

As commonly done, the values p, g, q are both system-wide and long-term. The same is true for the RSA public exponent e in the case of RSA-based key agreement: we assume that it is constant (e.g., 3 or $2^{16} + 1$) and known in advance. However, the device-specific DH key pairs (g^a, a) or the RSA composite modulo n and secret key d can be ephemeral (different for each key agreement). The DH public key g^a or the RSA modulo n are the authentication objects.

It is generally accepted that, for any b , factoring b -bit integers requires about the same amount of time as computing discrete logarithms in $(b - x)$ -bit fields, where x is a small constant, roughly around 20 [39]. Therefore, in our experimental and security analysis we use 1024-bit moduli for both RSA and DH.