# *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}(at)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 audio 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.*

## 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, wireless projector, network access point or another target 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 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 devices communicating 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 these 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. Section 4 describes the key elements of L&C. Sections 5 discusses unidirectional authentication using L&C. The L&C prototype implementation and its performance are discussed in Section 6 and 7, respectively. In [10], we list further related work, discuss the use of L&Cfor presence confirmation, and provide detailed analysis of our scheme's security properties.

## 2 Related Work

The well-known Diffie-Hellman key exchange protocol [24] 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 [2, 13, 15]. A number of enhancements to the Diffie-Hellman protocol have been developed. Bellovin and Merrit proposed the encrypted key exchange (EKE) protocol [19] to prevent MITM attacks. However, it requires both parties to possess a secret password *a priori*. Although many EKE refinements have been proposed [18, 21, 22], 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. For peer device authentication, Stajano and Anderson proposed a method for establishing keys by means of a link created through physical contact [6]. However, due to the diversity in personal devices, it is impractical to expect all devices to have suitable physical interfaces. Likewise, it may also be infeasible to lug around connection interfaces and or interface converters for various

---

[1]Although 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.

devices. Balfanz, et al. [5] extended this approach by using location-limited wireless infrared secondary channels. Since the human user is unable to directly verify which devices are communicating over the infrared channel, MITM attacks are still feasible by sophisticated attackers. Capkun et al. [20] proposed a further extension to allow two previously unassociated devices to establish a key utilizing one-hop transitive trust.

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) including [8, 9, 16] 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 [16] suggests) rendering their comparison a cumbersome and error-prone task for humans. [8] improves on [16] 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. More recently, Cagalj et al [4] tackled the problem of user-friendly mutual authentication with three commitment-based [23] techniques.

One notable recent result is the *Seeing-is-Believing* (SiB) system proposed by McCune, et al. [12]. 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.

## 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 four-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.
*4.* The use of camera-equipped devices is typically prohibited in high security facilities, such as military bases.

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 allow 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.

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 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 are complementary.

## 4  Main Elements of L&C

In this section, we describe the main elements of the Loud-and-Clear (L&C) system.

As a first application, we focus on the authentication of a

target device to a personal device, assuming no prior context between the two devices. We assume that in this, and most use scenarios, identification of the communicating devices is performed visually or tactilely 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.

## 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.

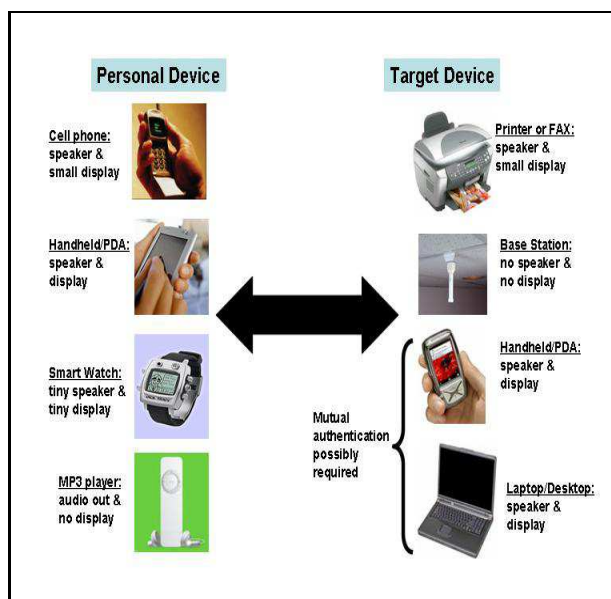Figure 1 depicts some anticipated L&C use scenarios.



**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 L&C use scenarios 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 [11] 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.

Table 1 shows the types of user requirements corresponding to possible personal-target device combinations.[2] 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.

|  | | Personal Device | | Target Device | |
|---|---|---|---|---|---|
| Row | Use Type | 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.

## 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 generated in L&C are based on the "MadLib" puzzles commonly used by children[3]. 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 or MD5 hash algorithms which are still believed to be second pre-image resistant. To achieve the desired level of entropy we take a prefix of the hash output.

Similar to the S/KEY One-Time Password System [17], the truncated hash is divided into 10-bit sections. The number of 10-bit sections becomes the final number of words contributing entropy to the MadLib sentence. (For example, using SHA-1 with 80 bits of entropy would result in a MadLib sentence containing 8 S/KEY words.)

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 FRAUDU-LENTLY CRUSH-ed over the CREEPY ARCTIC-TERN.

S/KEY generated words were not meant to be spoken, since there may be some similar-sounding words in the S/KEY word list. 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 [14], 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 is 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 [7] 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 converts bit sequences to auditorially robust words —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 proceed with describing some example use cases in detail.

---

[2]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.

[3]In 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.

# 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 Implementation

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* Java-based Programming System [3].[4] 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 the 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.

## 6.1 Implementation Details

Owing to its modular 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 Digalo [5], which is a simple lightweight clipboard reader that uses the *Elan Speech Engine*. [6] 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 [1] 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: *http://www.ics.uci.edu/ccsp/lac*.

---

[4]Ewe 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.

[5]See: *www.digalo.com*.
[6]See: *www.elantts.com*.

# 7  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 the L&C public key authentication completes, 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 $e$ 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 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 [10]).

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"); (2) pressing "Enter" and "Connect" buttons; (3) aligning the devices if the Infrared channel is used; and (4) 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 |
|---|---|

Note: all entries below refer to Table 2.

(A) Type 1

| Row 6 Col. 3 | 5.6 |
|---|---|
| Row 8 Col. 3 | 365.6 |
| Row 9 Col. 4 | 4,791 |
| Row 9 Col. 3 | 4,637 |
| TOTAL TIME | 32,044 |

(B) Type 2

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

(C) Type 4

| Row 6 Col. 3 | 5.6 |
|---|---|
| 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 per-

forming 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. It involves unidirectional 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 unidirectional 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 displayed MadLib sentences. We did not measure TYPE 3, due to its similarity to TYPE 2.

## 7.1  Performance Analysis

Table 2 illustrates that the overall cost is dominated by the Ewe VM and GUI initializations and the MadLib vocalizations. The initializations can be omitted if multiple L&C sessions take place after a single initialization or if L&C is pre-initialized. Since the time to speak out a MadLib sentence is proportional to number of syllables in each S/KEY-generated word, 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 non-user-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.

Our experiments suggest that the most plausible way to reduce protocol overhead is to shorten MadLib sentences and speed up the TTS engine initialization time. We conclude that L&C is a viable solution on platforms with moderate computation and communication capabilities.

## 7.2  Performance Improvements

Below, we focus on improving the total time of a L&C session by shortening the length of the MadLib sentences. A US patent addressing a similar problem to ours [16] proposed that upon completion of a Diffie-Hellman exchange, both devices one-way hash the agreed-upon session key. Then they truncate the hash to the desired length by taking its $t$ 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. That proposal used $t = 16, 32$ which corresponds to 2 to 4 S/KEY words. However, using a hash-code that short results to a serious security weakness.[7] Based on the assumed adversarial capabilities (see [10]), $t = 50$ and $t = 80$ provide the necessary security when ephemeral public keys and one-year-term keys are used, respectively.[8]

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.

A technique proposed in [8] prevents attackers from finding second pre-images for long-term public keys, without requiring the users to examine 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 and presents the MadLib sentence for the check-value (five words) and a 4 digit hexadecimal sequence for the random key.

*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 presents 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

---

[7]The 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 $b'$ such that $h_t(g^{a'b}) = h_t(g^{ab'})$, 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.

[8]The seven S/KEY word generated MadLibs used in the evaluation provide sufficient security if the public keys are renewed daily.

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, 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 [4], 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 [10]).

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 a 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 Conclusions

This paper discussed 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 a prototype implementation.

# References

[1] Bouncy Castle Crypto APIs. http://www.bouncycastle.org/index.html.

[2] A. J. Menezes, P. C. van Oorschot, and S. A. Vanstone. *Handbook of Applied Cryptography*. CRC Press Series on Discrete Mathematics and its Applications. CRC Press, 1997.

[3] M. Brereton. Ewe java vm for pocketpc. http://www.ewesoft.com/.

[4] M. Cagalj, S. Capkun, and J. Hubaux. Key agreement in peer-to-peer wireless networks. IEEE Special Issues on Cryptography and Security), 2006.

[5] D. Balfanz, D.K. Smetters, P. 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.

[6] F. Stajano and R. Anderson. The resurrecting duckling: Security issues for ad-hoc wireless networks. In *Security Protocols, 7th International Workshop*, 1999.

[7] M. R. Garey and D. S. Johnson. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., New York, NY, USA, 1990.

[8] C. Gehrmann, C. Mitchell, and K. Nyberg. Manual authentication for wireless devices. RSA Cryptobytes, Vol. 7, No. 1, pp. 2937, 2004.

[9] C. Gehrmann and K. Nyberg. Security in personal area networks. In C. J. Mitchell, editor, Security for Mobility, pages 191230. IEE, London, 2004.

[10] M. T. Goodrich, M. Sirivianos, J. Solis, G. Tsudik, and E. Uzun. Loud and clear: Human-verifiable authentication based on audio. Cryptology ePrint Archive, Report 2005/428, 2005. http://eprint.iacr.org.

[11] D. L. Greg Kochanski and C. Shih. A reverse turing test using speech. In *Seventh International Conference on Spoken Language*, September 2002. Denver, Colorado.

[12] J. McCune, A. Perrig, M. K. Reiter. Seeing-Is-Believing: Using Camera Phones for Human-Verifiable Authentication. In *IEEE Symposium on Security and Privacy*, pages pp. 110–124, 2005.

[13] M. Jakobsson and S. Wetzel. Security weaknesses in bluetooth. *Lecture Notes in Computer Science*, 2020:176+, 2001.

[14] P. Juola and P. Zimmermann. Whole-word phonetic distances and the pgpfone alphabet. In *Fourth International Conference on Spoken Language Processing (ICSLP)*, volume 1, 1996.

[15] D. Kugler. Man in the middle attacks on bluetooth. Financial Cryptography '03, Long Beach, 2003. Lecture Notes in Computer Science, Springer-Verlag.

[16] D. P. Maher. Secure communication method and apparatus. U.S. Patent Number 5,450,493 September 1995.

[17] N. Haller. Rfc1760: The s/key one-time password system, 1995.

[18] P. MacKenzie, S. Patel, and R. Swaminathan. Password authenticated key exchange based on RSA. In *ASIACRYPT*, pages 599–613, 2000.

[19] S. Bellovin and M. 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*, pages 244–250, 1993.

[20] S. Capkun, J. Hubaux, and L. Buttyan. Mobility helps security in ad hoc networks. In *ACM MobiHoc 2003*, June 2003.

[21] T. Wu. The secure remote password protocol. In *Network and Distributed System Security Symposium*, February 1999.

[22] V. Boyko, P. MacKenzie, and S. Patel. Provably secure password authentication and key exchange using diffie-hellman. volume 1807 of Lecture Notes in Computer Science, pages 156–171, 2000.

[23] S. Vaudenay. Secure communications over insecure channels based on short authenticated strings. In *CRYPTO*, pages 309–326, 2005.

[24] W. Diffie and M. E. Hellman. New directions in cryptography. In *IEEE Transactions on Information Theory*, pages IT–22(6):644–654, 1976.

COMPUTER SOCIETY