

# Voice Biometrics over the Internet in the Framework of COST Action 275

**Laurent Besacier,<sup>1</sup> Aladdin M. Ariyaeenia,<sup>2</sup> John S. Mason,<sup>3</sup> Jean-François Bonastre,<sup>4</sup>  
Pedro Mayorga,<sup>1</sup> Corinne Fredouille,<sup>4</sup> Sylvain Meignier,<sup>4</sup> Johann Siau,<sup>2</sup>  
Nicholas W. D. Evans,<sup>5</sup> Roland Auckenthaler,<sup>5</sup> and Robert Stapert<sup>6</sup>**

<sup>1</sup> CLIPS/IMAG, 38041 Grenoble Cedex 9, France

Emails: laurent.besacier@imag.fr; pedro.mayorga-ortiz@imag.fr

<sup>2</sup> Department of Electronic, Communication and Electrical engineering, University of Hertfordshire, Hatfield, AL10 9AB, UK

Emails: a.m.ariyaeenia@herts.ac.uk; j.siau@herts.ac.uk

<sup>3</sup> Department of Electrical and Electronic Engineering, University of Wales Swansea, Swansea SA2 8PP, UK

Email: j.s.d.mason@swansea.ac.uk

<sup>4</sup> LIA, University of Avignon, 84911 Avignon Cedex 9, France

Emails: jean.francois.bonastre@lia.univ-avignon.fr; corinne.fredouille@lia.univ-avignon.fr;  
sylvain.meignier@lia.univ-avignon.fr

<sup>5</sup> School of Engineering, University of Wales Swansea, Swansea SA2 8PP, UK

Emails: n.w.d.evans@swan.ac.uk; eeaucken@swansea.ac.uk

<sup>6</sup> Aculab, Milton Keynes, MK1 1PT, UK

Email: robert.stapert@aculab.com

Received 1 December 2002; Revised 3 September 2003

The emerging field of biometric authentication over the Internet requires both robust person authentication and secure computer network protocols. This paper presents investigations of vocal biometric person authentication over the Internet, both at the protocol and authentication robustness levels. As part of this study, an appropriate client-server architecture for biometrics on the Internet is proposed and implemented. It is shown that the transmission of raw biometric data in this application is likely to result in unacceptably long delays in the process. On the other hand, by using data models (or features), the transmission time can be reduced to an acceptable level. The use of encryption/decryption for enhancing the data security in the proposed client-server link and its effects on the transmission time are also examined. Furthermore, the scope of the investigations includes an analysis of the effects of packet loss and speech coding on speaker verification performance. It is experimentally demonstrated that whilst the adverse effects of packet loss can be negligible, the encoding of speech, particularly at a low bit rate, can reduce the verification accuracy considerably. The paper details the experimental investigations conducted and presents an analysis of the results.

**Keywords and phrases:** voice biometrics, speaker verification, packet loss, compression, Internet.

## 1. INTRODUCTION

The ever-increasing use of the Internet-enabled devices is resulting in normal activities in day-to-day life, such as banking and shopping, being conducted without face-to-face or personal contacts. A natural consequence of this is the obsolescence of certain conventional means of identification. Examples of these are photo ID cards and passports. On the other hand, the conventional authentication means such as personal identification numbers and passwords, which are equally applicable to local and remote identity verification, can be easily compromised or forgotten. In view of the above,

it appears that biometrics is the only means that can satisfy the requirements for remote identity verification in terms of both appropriateness and reliability. This is because firstly, biometric data can be easily captured, stored, processed, and described electronically. Secondly, it uses an intrinsic aspect of a human being for identity verification. Consequently, it is not so susceptible to fraud as passwords or personal identification numbers.

The deployment of biometrics on the Internet, however, is a multidisciplinary task. It involves person authentication techniques based on signal processing, statistical modelling, and mathematical fusion methods, as well as data

communications, computer networks, communication protocols, and online data security.

The necessity for the latter discipline is due to the fact that an online robust biometric authentication strategy would be of little or no value if, for instance, hackers could break into the personal identification server to control the verification of their pretended identities, or could access personal identification data transmitted over the network.

The original aim of the Internet was to provide a means of sharing information, thus security was not of major concern. As the Internet has evolved, many security implications and bandwidth issues have arisen. There are many potential threats to any system that relies on the Internet as a communication medium. The potential benefits of biometric identity verification over the Internet have highlighted issues of security and network performance that need to be tackled more effectively [1].

In general, network performance varies widely with the geographical location of the clients, server type, and network resources. There is variation in the response time from session to session even if the connection is made to the same server. This is because in each session, data packets may travel through a different route [2]. There is a difference in the performance of the dial-up Internet service, integrated subscriber digital network (ISDN), asymmetric digital subscriber line (ADSL), cable modem, and leased line as they all have a different bandwidth and response time. This will undoubtedly affect the performance of biometric verification systems in terms of speed, reliability, and the quality of service.

Over IP networks, both speech and image-based biometrics are viable alternative approaches to verification. Focusing on speech biometrics, some predictions for the year 2005 show that 10% of voice traffic will be over IP. This means that speaker verification technology will have to face new problems. The most common architecture seems to be client-server-based where a distant speaker verification server is remotely accessed by the client for authentication. In this scenario, the speech signal is transmitted from the client terminal to a remote speaker verification server. Coding of the speech signal is then generally necessary to reduce transmission delays and to respect bandwidth constraints. Many problems can appear with this kind of architecture, particularly when the transmission is made via the Internet:

- (i) firstly, transcoding (the process of coding and decoding) modifies the spectral characteristics of the speech signal, and thereby can adversely affect the speaker verification performance;
- (ii) secondly, transmission errors can occur on the transmission line: thus, data packets can be lost (e.g., with UDP transport protocols which do not implement any error recovery);
- (iii) thirdly, the time response of the system is increased by coding, transmission, and possible error recovery processes. This delay (termed "jitter" as used in the domain of computer networks) can be potentially very disturbing. For example, in some applications (e.g.,

man-machine dialogue), speaker verification is only one subsystem amongst a number of other subsystems. In such cases, the effective operation of the whole system depends heavily on the response time of the individual subsystems;

- (iv) finally, speech packets (or other personal information) transmitted over IP could be intercepted and captured by impostors, and subsequently used, for instance, for fraudulent access authorisation.

To our knowledge, this paper is the first to present an overview of issues and problems in the above area. These include architecture and protocol considerations (Section 2), speaker verification robustness to speech coding and packet loss over IP networks (Section 3), and wireless mobile devices (Section 4). This work is currently conducted in the framework of COST Action 275 (<http://www.fub.it/cost275/>).

## 2. ARCHITECTURE AND PROTOCOL CONSIDERATIONS IN BIOMETRICS OVER THE INTERNET

This part details an analysis carried out to determine the right balance in the transmission method for the purpose of implementing applications involving biometric verification. These tests were conducted in different geographical locations within the UK. However, most of the local area network (LAN) tests were carried out in the premises of the University of Hertfordshire.

### 2.1. Biometrics applied

The raw biometric data can have different sizes depending on its type. For instance, voice or face biometric datasets are considerably larger than that of fingerprint. In any case, the data contains the identity of an individual and should be treated with utmost care. Therefore, it is necessary to have an appropriate architecture and method of transmission in order to provide a high level of protection against uncertainties.

#### 2.1.1. Client-server architecture

An effective client-server structure for biometrics on the Internet has recently been proposed by some authors of this paper [3]. This realisation (Figure 1) consists of 3 distinct components, each performing a specific task. The client part consists of users (clients) requesting appropriate services from the server. A main role of the server is to respond to these requests. However, from time to time, it itself becomes a client to the central database and requests services from it.

The modular nature of the proposed structure is also necessary for performing software updating effectively. For example, the client module dynamically obtains information relevant to its process, and the updates to its software are provided by the server. As a result, it is ensured that the client software will always be up-to-date, and modifications or improvements can be gradually rolled in.

In order to maintain data integrity, the transmission channel needs to be secured and encrypted. This will ensure

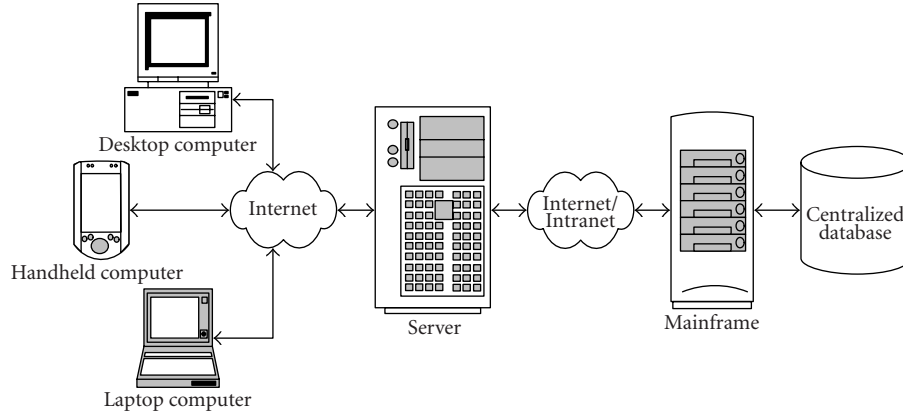
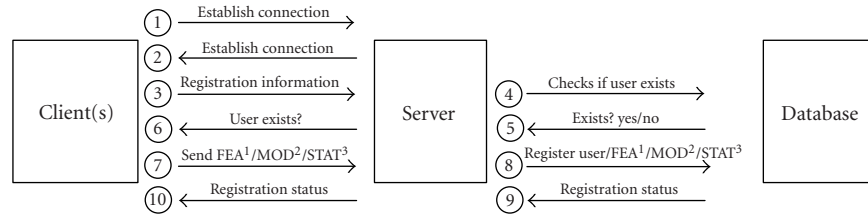
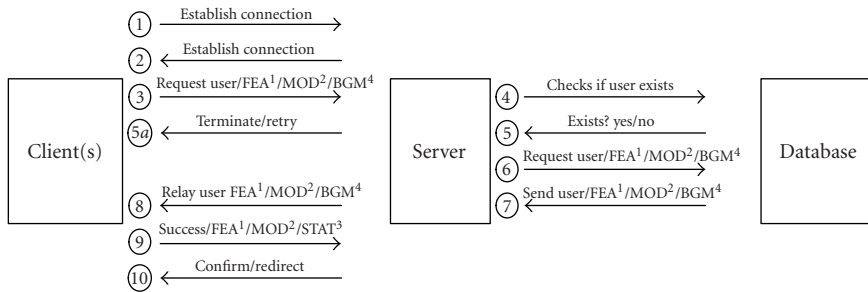


FIGURE 1: Client-server architecture.



FEA<sup>1</sup> (features)  
 MOD<sup>2</sup> (models)  
 STAT<sup>3</sup> (statistics/scores)

(a)



FEA<sup>1</sup> (features)  
 MOD<sup>2</sup> (models)  
 STAT<sup>3</sup> (statistics/scores)  
 BGM<sup>4</sup> (background model)

(b)

FIGURE 2: Proposed client-server architecture. (a) Enrolment process. (b) Verification process.

that data sent from the client to the server and vice versa will be of no use to others even if they breach the system.

Figure 2 illustrates the operation of the proposed system in terms of its enrollment and verification processes. It should be noted that although the system is ideally suited to speaker verification, it could also be adapted to suit other types of biometrics. The operation can be described as follows.

The database acts as the central storage area for all biometric data and also as a server to the main server. Each server has its unique identifier that allows its connection to the database. All communications between the server and database are secured and encrypted. Distributed/different servers from different geographical locations can therefore connect to the central database through a fast network link.

During the enrollment process, the client initially establishes a connection with the server. This is known as the handshaking process in which the client and server establish the identity of both machines for that particular session. The encryption key (Section 2.1.3) is also exchanged at this time. The registration information is then sent to the server. Once a confirmation is obtained from the server that the user does not exist in the system, the client is prompted to send the biometric features, models, and statistics over to the server to be enrolled. These are encrypted before transmission. The server then forwards this information to the database and thus enrolling the user to the system.

When a user returns to verify his/her identity, the client machine establishes a connection with the server, whereby during the handshaking process, a different key will be allocated to secure the connection for the session. The client then requests the server to provide data files associated with the user. The server then requests the relevant information from the central database and relays the data back to the client. The client machine uses this information to perform a verification test. If the test result is positive, the statistics regarding the success of the verification is sent back to the server to be stored into the central database.

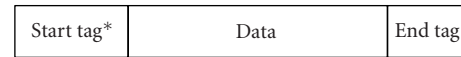
Depending on the level of security required, the function of the client machine, and the location of the client machine, some operations can be adapted to optimise the performance-to-security ratio appropriately. For example, when a home PC is used, the data files can be stored on the local computer for later use. This will result in reducing the amount of data transfer necessary between the client and the server. However, when the client uses a station which is not registered as his/her own, then the data files provided by the server will need to be removed from the client station after each process is completed in order to improve the security measures.

An advantage of the above architecture is that it will allow, and accommodate, future expandability and upgradeability beyond that achievable with a conventional software-based system architecture. Additionally, unlike some newly developed online recognition systems (<http://www.biometrika.it>), the proposed architecture eliminates the need for the installation of software on local terminals. This enhances the usability of the online recognition system considerably as it allows access from any station and any location.

Moreover, the proposed architecture requires only minimal data to be transmitted between client-server-database, as opposed to the transmission of the full raw biometric data. The emergence of load-balancing and distributed systems technology provides the possibility of having servers distributed at different remote locations. This in turn further reduces the time-lag in client-server communications.

### 2.1.2. Data format

As in most client-server architectures, a set of instructions is needed to enable communications between the client software and the server software. The instructions for the system follow a format similar to that shown in Figure 3. The start



\*Start tag contains either control, data, or key tags

FIGURE 3: Data format tags.

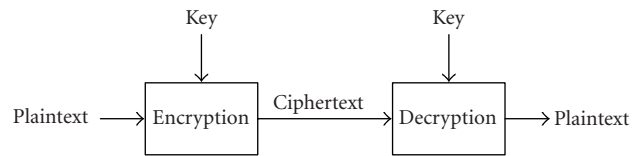


FIGURE 4: Encryption/decryption process.

tag contains one of control, data, or key tags as appropriate for the correct operation of the system.

It is worth noting that the biometric information transferred should be in the form of characteristic features rather than raw data. This will reduce the size of the data to be transferred. Moreover, with this approach, the load on the server can be reduced by performing parts of the processing on the client machine.

### 2.1.3. Data security

The transmission of data over the network requires some form of security measure. Sensitive data such as biometrics needs to be encrypted to prevent others from misusing it. Therefore, the link between the client and server has to be secure throughout the entire process to prevent access or attacks from a hostile source.

To secure the link between the client and the server effectively, the data transmitted between them needs to be in encrypted form. Encryption is a process of disguising/ciphering a message which hides its contents by representing it in a different form. For the purpose of decryption, the exact key used for the encryption process will be needed to restore the original message. Without knowing the key, it will be practically impossible to access the message contents. This process is summarized in Figure 4.

A well-known algorithm for encrypting and decrypting messages is Blowfish [4]. This algorithm is in the public domain and is considered for the purpose of this study. A main advantage of Blowfish is that it is significantly faster than data encryption standard (DES) [5]. A description of Blowfish is presented in the following section.

### 2.1.4. Blowfish

Blowfish is a 64-bit block cipher, and the algorithm consists of two parts. These are a key-expansion part and a data-encryption part. Key expansion converts a key of at most 448 bits into several subkey arrays in a total of 4168 bytes. The data is then encrypted via a 16-round Feistel network, where each round consists of a key-dependent permutation and a key- and data-dependent substitution. All operations are XORs and additions on 32-bit words. The only

TABLE 1: Dependence of the transmission time(s) on the file size and connection type.

File size (bytes)	Connection					
	Dial-up 56 k	Cable/DSL 512 k	Cable/DSL 1 M	LAN 10 M	LAN 100 M	LAN 1 G
87 k	12.43	1.36	0.68	0.07	0.01	$0.6 \times 10^{-3}$
130 k	18.57	2.03	1.02	0.10	0.01	$1.0 \times 10^{-3}$
173 k	24.71	2.70	1.35	0.14	0.01	$1.4 \times 10^{-3}$
216 k	30.86	3.38	1.69	0.17	0.02	$1.7 \times 10^{-3}$
259 k	37.00	4.05	2.02	0.20	0.02	$2.0 \times 10^{-3}$
302 k	43.14	4.72	2.36	0.24	0.02	$2.4 \times 10^{-3}$
345 k	49.29	5.39	2.70	0.27	0.03	$2.7 \times 10^{-3}$
388 k	55.43	6.06	3.03	0.30	0.03	$3.0 \times 10^{-3}$
431 k	61.57	6.73	3.37	0.34	0.03	$3.4 \times 10^{-3}$
517 k	73.86	8.08	4.04	0.40	0.04	$4.0 \times 10^{-3}$
603 k	86.14	9.42	4.71	0.47	0.05	$4.7 \times 10^{-3}$
690 k	98.57	10.78	5.39	0.54	0.05	$5.4 \times 10^{-3}$
776 k	110.86	12.13	6.06	0.61	0.06	$6.1 \times 10^{-3}$
862 k	123.14	13.47	6.73	0.67	0.07	$6.7 \times 10^{-3}$
1024 k	146.29	16.00	8.00	0.80	0.08	$8.0 \times 10^{-3}$

additional operations are four indexed array data lookups per round.

Blowfish uses a large number of subkeys for encryption or decryption and these keys must be precomputed before any of the above processes can be carried out. The generation of the subkeys involves two arrays consisting of eighteen 32-bit  $P$ -arrays subkeys  $P_1 \cdots P_{18}$  and four 32-bit  $S$ -boxes with 256 entries each.

The calculation of the subkeys is detailed in Schneier's paper [4]. In general, generating the subkeys is a computationally expensive process and requires a total of 521 iterations. However, these keys can then be stored and reused.

## 2.2. Experimental analysis

The most common connection to the Internet is normally via a dial-up service which ideally offers a maximum transmission speed of 56 kbps. However, cable/ADSL services are becoming more and more available. In an ideal situation, these offer services with transmission speeds of up to 1 Mbps downstream (receiving data) and 512 kbps upstream (sending data). However, the most common transmission speeds of these for receiving and sending data are 512 kbps and 256 kbps, respectively. It should also be noted that these transmission rates might vary considerably during a given connection.

### 2.2.1. Theoretical transmission rates

The basic approach to calculate the time taken to transmit a file from one location to another via the Internet is based on the following equation:

$$T_s = \frac{Fsz \times 8}{Cnx}, \quad (1)$$

where  $T_s$  is the time taken in seconds,  $Fsz$  is the file size in bytes, and  $Cnx$  is the connection speed in bps.

The above equation assumes an ideal situation where the connection to the Internet and to the destination servers is achieved at the maximum throughput. This, however, is not the actual case on a day-to-day basis.

A comparison of the calculated theoretical transmission time for different file sizes and different connection types is presented in Table 1.

As observed in this table, even in an ideal situation, the use of a dial-up connection involves relatively a long transmission time.

### 2.2.2. Experimental transmission rates

Experiments were conducted at different times using two types of common Internet connections with the file size varying from 4 kb to 900 kb. The files used were signals generated from white noise. These audio files were of 1 to 10 seconds in length. The two types of connection used were a 56 k dial-up connection service and a LAN. The results of this experimental study are given in Figure 5. As it is observed, the transmission time in practice is significantly longer than that suggested theoretically.

The results in Figure 5 clearly indicate that verification over the Internet is unfavourably influenced by the performance of the network. To minimize this, it seems advantageous to compress data before its transmission.

The next set of experiments was based on the transmission of audio models rather than raw data. The previous set of white noise files (Section 2.2.2) was preprocessed and the features were extracted using LPCC-12. These were used to generate audio models based on a VQ with a codebook size of 64. The results of this study are presented in Table 2. As observed, due to the use of VQ, considerable reduction in the file size is achieved. This in turn has resulted in significant reduction in transmission time.



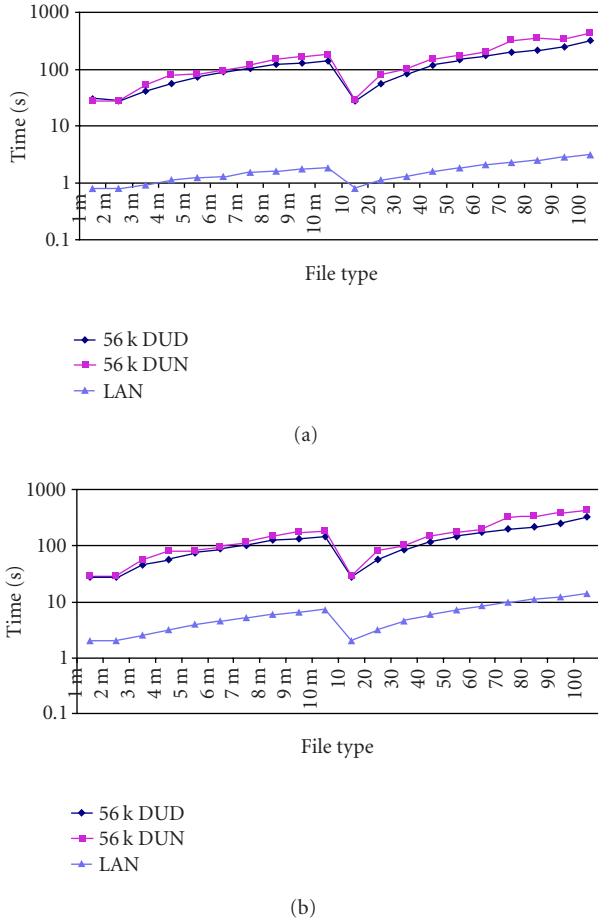


FIGURE 5: Experimental transmission rates (DUD: dial-up daytime; DUN: dial-up nighttime). (a) Transmission times without encryption. (b) Transmission times with encryption.

As part of this study, a second set of experiments was conducted based on the encryption of VQ files using the Blowfish algorithm. The results of this investigation are also shown in Table 2. It is seen that there is a slight increase in the overall transmission time in this case. This is due to the initial processing time needed to prepare the data prior to transmission and the time taken to decrypt the data at the receiver. The resultant increase in the overall transmission time is negligible and often not noticeable.

These experimental results indicate the difficulties introduced by the transmission of raw data over the Internet, especially when the file sizes are too large. The results presented were based on the use of audio signal files. It should be noted that image-based biometric data files are of considerably larger sizes. The transmission of such raw files over the Internet may sometimes result in unacceptably long delays in the verification process.

### 2.3. Comments

A client-server architecture for biometric verification over the Internet has been proposed and described in detail. Based

TABLE 2: Transmission time for 4 KB audio models (DUD: dial-up daytime; DUN: dial-up nighttime).

LPCC12 VQ64	Transmission time(s)	
	Without encryption	With encryption
56 k DUD	1.9	2.3
56 k DUN	2.6	2.7
LAN	0.1	0.2

on an analysis of the characteristics of the proposed architecture, its advantages have been discussed, and it has been shown that it provides a practical and systematic approach to the implementation of biometric verification on the Internet. Using a set of experimental investigations, it has been shown that, in practice, it may not be feasible to transmit raw biometric data over the Internet as this can cause unacceptably long delays in the process. It has been demonstrated that the transmission of data models (or features) instead of raw material will significantly reduce the transmission time. Another possibility is to compress biometric data before its transmission. Such compression, however, may unfavourably influence the robustness of biometric techniques (see the next part). Finally, it has been argued that the client-server link should be made secure by encrypting the data before its transmission. It has been shown that the increase in the overall transmission time due to this process is relatively small.

## 3. SPEAKER VERIFICATION EXPERIMENTS OVER IP NETWORKS

In Section 2, it has been notably shown that transmitting raw biometric data over the Internet may lead to unacceptably long delays. However, recently, considerable progress has been achieved in transmitting voice over the Internet for communication purposes. Thus, this section proposes a methodology for evaluating the speaker verification performance over IP network. The idea is to duplicate an existing and well-known database used for speaker verification (XM2VTS) by passing its speech signals through different coders and different network conditions representative of what can occur over the Internet. Some partners of COST 275 are also evaluating the influence of image and video compression on face recognition performance, again using XM2VTS as it is a multimodal database. Section 3.1 is dedicated to the database description and to the degradation methodology adopted, whereas Section 3.2 presents the speaker verification system and some results obtained with this IP-degraded version of XM2VTS.

### 3.1. Database used and degradation methodology

#### 3.1.1. XM2VTS database

In acquiring the XM2VTS database (<http://www.ee.surrey.ac.uk/Research/VSSP/xm2vtsdb/>), 295 volunteers from the University of Surrey visited a recording studio four times at approximately one-month intervals. On each visit, (session)

two recordings (shots) were made. The first shot consisted of speech while the second consisted of rotating head movements. Digital video equipment was used to capture the entire database. At the third session, a high-precision 3D model of the subjects head was also built using an active stereo system provided by the Turing Institute. We have chosen this database since many partners of COST Action 275 already use it. The work described in this paper was made on its speech part, where the subjects were asked to read three sentences twice. The three sentences remained the same throughout all four recording sessions and a total of 7080 speech files were made available on 4 CD-ROMs. The audio, which had originally been stored in mono, 16 bit, 32 kHz PCM wave files, was down-sampled to 8 kHz. This is the input sampling frequency required in the speech codecs considered in this study.

### 3.1.2. Codec used

H323 is a standard for transmitting voice and video. A famous H323 videoconferencing software is for example NetMeeting™. H323 is commonly used to transmit video and voice over IP networks. The audio codecs used in this standard are G711, G722, G723.1, G728, and G729. We propose to use in our experiments the codec which has the lowest bit rate: G723.1 (6.4 and 5.3 kbps), and the one with the highest bit rate: G711 (64 kbps: 8 kHz, 8 bits). Influence of these codecs on speech recognition was evaluated in a former study we made [6], it is thus very exciting to know what will be the results on the speaker verification task.

### 3.1.3. Packet loss

#### *Simulation with the Gilbert model*

There are two main transport protocols used on IP networks. These are UDP and TCP. While UDP protocol does not allow any recovery of transmission errors, TCP includes some error recovery processes. However, the transmission of speech via TCP connections is not very realistic. This is due to the requirement for real-time (or near real-time) operations in most speech-related applications [7]. As a result, the choice is limited to the use of UDP which involves packet loss problems. The process of audio packet loss can be simply characterised using a Gilbert model [8, 9] consisting of two states (Figure 6). One of the states (state 1) represents a packet loss and the other state (state 0) represents the case where packets are correctly transmitted. The transition probabilities in this statistical mode, as shown in Figure 6, are represented by  $p$  and  $q$ . In other words,  $p$  is the probability of going from state 0 to state 1 and  $q$  is the probability of going from state 1 to state 0.

Different values of  $p$  and  $q$  define different packet loss conditions that can occur on the Internet. The probability that  $n$  consecutive packets are lost is given by  $p(1 - q)^{n-1}$ . If  $(1 - q) > p$ , then the probability of losing a packet in state 1 (after having already lost a packet) is greater than the probability of losing a packet in state 0 (after having successfully received a packet) [9]. This is generally the case in data transmission on the Internet where packet losses occur

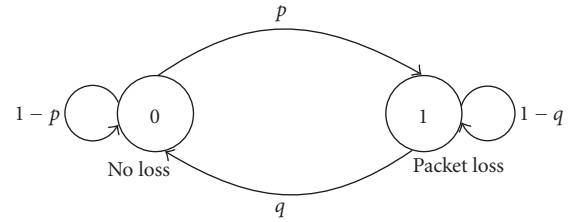


FIGURE 6: Gilbert model.

as bursts. Note that  $p + q$  is not necessarily equal to 1. When  $p$  and  $q$  parameters are fixed, the mean number of consecutive packets lost can be easily calculated as  $p/q^2$ . Of course, the larger this mean is, the more severe the degradation is. Different values of  $p$  and  $q$  representing different network conditions considered in this study are presented in Table 3 [8, 9].

#### *Real-conditions packet loss*

In order to investigate the effects of real network conditions as well, it was decided to play and record the whole speech part of XM2VTS through the network. This was carried out by playing the speech dataset into a computer which was set up for videoconferencing. For this purpose, a transatlantic connection was established between France and Mexico using videoconferencing software. The microphone on the French site was then replaced with the audio output of a computer playing the speech material in XM2VTS. Due to numerous network breakdowns, the transmission of material had to be conducted using several different connections established on different days and at different times. This, of course, provided variations in network conditions that occur in the case of real applications. Table 3 presents a summary of the different coders and simulated network conditions that were considered.

- (i) Two degraded versions of XM2VTS were obtained by applying G711 and G723.1 codecs alone without any packet loss.
- (ii) Six degraded versions of XM2VTS were obtained using simulated packet loss conditions: 2 conditions (average/bad)  $\times$  3 speech qualities (clean/G711/G723.1). The simulated average and bad network conditions considered in this study corresponded to 9% and 30% speech packet loss rates, respectively. Each packet contained 30 milliseconds of speech which was consistent with the duration proposed in Real Time Protocol (RTP) (used under H323).
- (iii) One degraded version of XM2VTS based on real network conditions. The transmission was spread from 12/9/02 to 1/10/02 and the mean packet loss rate was 15%. The detailed packet loss conditions for each part of the database are described in Figure 7. Each bar corresponds to a different transmission day and thus to a different transmission condition. We see that in the worst cases, real packet loss rate is around 30%; this

TABLE 3: Summary of the simulated IP degradation plan (3 codecs \* 3 network conditions give 9 different degradations).

Codecs	None (128 kbps)	G711 (64 kbps)	G723.1 (5.3 kbps)
Network condition	No packet loss	Average $p = 0.1; q = 0.7$	Bad $p = 0.25; q = 0.4$

figure corresponds approximately to the mean packet loss rate measured after simulated IP degradation with  $p = 0.25$  and  $q = 0.4$  (called bad condition in Table 3). On the other hand, in the best cases, real packet loss rate is around 10% and even less; this corresponds approximately to our simulated “average” condition ( $p = 0.1; q = 0.7$  in Table 3) for which mean packet loss rate is around 9%.

### 3.2. Speaker verification experiments with the ELISA system

The ELISA consortium groups several public laboratories working on speaker recognition. One of the main objectives of the consortium is to emphasize assessment of performance. Particularly, the consortium has developed a common speaker verification system which has been used for participating at various NIST speaker verification evaluations campaigns [10, 11].

ELISA system is a complete framework designed for speaker verification. It is a Gaussian mixture model (GMM) based system [12] including audio parameterisation as well as score normalization techniques for speaker verification.

This system was presented at NIST from 1998 to 2002 and showed the state-of-the-art performance. ELISA is now collaborating with COST Action 275 concerning performance assessment of multimodal person authentication systems over the Internet. ELISA evaluated the speaker verification performance using the COST 275 dedicated database detailed in Section 3.1.

#### 3.2.1. Speaker verification protocol on XM2VTS

For the purpose of this investigation, the Lausanne protocol (configuration 2) is adopted. This has already been defined for the XM2VTS database. There are 199 clients in the XM2VTS database. The training of the client models is carried out using full session 1 and full session 2 of the client part of XM2VTS. Test accesses of 398 clients are obtained using full session 4 ( $\times 2$  shots) of the client part. Using the impostor part of the database (70 impostors  $\times 4$  sessions  $\times 2$  shots  $\times 199$  clients = 111440 impostor accesses) 111440 impostor accesses are obtained. The 25 evaluation impostors of XM2VTS are used to develop a world model. The text-independent speaker verification experiments are conducted in matched conditions (same training/test conditions).

#### 3.2.2. ELISA system on XM2VTS

The ELISA system on XM2VTS is based on the LIA system presented to NIST 2002 speaker recognition evaluation. The speaker verification system uses 32 parameters: 16 linear fre-

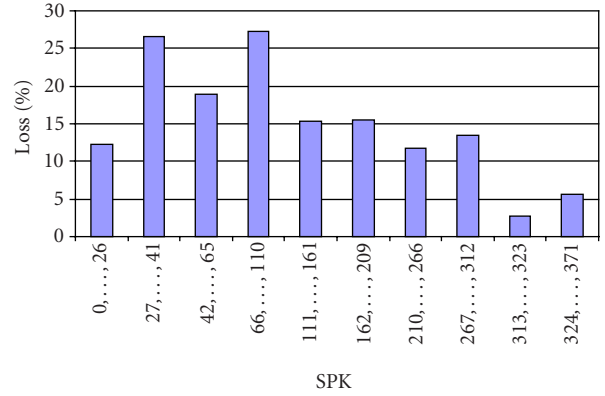


FIGURE 7: Packet loss measurements for real transmission over IP (different groups of speakers SPK represent different connections).

quency cepstral coefficients (LFCC) + 16 DeltaLFCC. Silence frame removal is applied before centring (CMS) and reducing vectors.

For the world model, 128-Gaussian component GMM was trained using Switchboard II phase II data (8 kHz land-line telephone) and then adapted (MAP [13], mean only) on XM2VTS data (25 evaluation impostors set). The client models are 128-Gaussian component GMM developed by adapting (MAP, mean only) the previous world model.

Decision logic is based on using the conventional log likelihood ratio (LLR). No LLR normalisation such as  $Z_{\text{norm}}$  [14],  $T_{\text{norm}}$  [15], or  $D_{\text{norm}}$  [16] is applied before the decision process.

#### 3.2.3. Results

The speaker verification performance with the simulated degraded versions of XM2VTS is presented in Table 4. We can see that whatever the packet loss level is (no packet loss, average condition, or bad condition), the equal error rate (EER) remains very low for clean speech (no codec) or slightly compressed speech (G711). Based on these results, it can be concluded that, even at a high rate, packet loss alone is not a significant problem for text-independent speaker verification. Comparing these results with those for speech recognition [17], it can be said that the speaker verification performance is far less sensitive to packet loss. On the other hand, the last column of Table 4 shows that the speaker verification performance is adversely affected when the speech material is encoded at low bit rates (e.g., using G723.1). In that case, packet loss increases the degradation. These results are in agreement with those in Section 4 of this paper, describing the performance of speaker verification over wireless mobile devices.



TABLE 4: Results (EER%) of the experiments using degraded XM2VTS.

Network condition	Codecs		
	Clean (128 kbps)	G711 (64 kbps)	G723.1 (5.3 kbps)
No packet loss	0.25%	0.25%	2.68%
Average Network condition $p = 0.1; q = 0.7$	0.25%	0.25%	6.28%
Bad Network condition $p = 0.25; q = 0.4$	0.50%	0.75%	9%

#### 4. SPEAKER VERIFICATION EXPERIMENTS OVER WIRELESS MOBILE DEVICES

Most wireless mobile networks are susceptible to packet loss to some degree. Whilst there exist many strategies to combat packet loss, such as retransmission or packet recovery [17, 18, 19], online identity verification applications may still operate effectively from semi real-time voice streams. This is possible because there is no intrinsic requirement on latency in the case of retransmission. In this part, speaker verification accuracy is assessed against the level of packet loss in wireless mobile devices.

The packet loss scenario is contrasted with degradation coming from additive noise. The degrading effect of ambient noise on automatic speech and speaker recognitions is widely acknowledged and known to be large even for relatively low noise levels. Thus a comparison is made between the two forms of degradation by using otherwise identical experimental conditions.

The remainder of this part is organised as follows. Section 4.1 addresses packet loss in typical wireless and IP networks and its effects on speaker verification. Section 4.2 addresses additive noise and speech enhancement.

Experimental work on the 2000-speaker SpeechDat Welsh [20] database is presented in Section 4.3 with results of experiments using both simulated packet loss and speech enhancement after contamination by additive real car noise.

##### 4.1. Packet loss in mobile networks

Some degree of packet loss is inherent in mobile networks. Lost packets might be caused by variable transmission conditions, or the hand-over between neighbouring cells as a wireless mobile device roams about the network.

Approaches dealing with packet loss recovery are generally controlled by the routing protocol adopted in the network architecture. For automatic speech recognition applications where time-sequence information is more critical, packet loss might have a significant impact on performance.

Lost packets might then be retransmitted or some form of compensation employed [17, 18, 19]. In contrast, as seen in Section 3, for speaker verification, a limited degree of

packet loss might not have a too detrimental effect, particularly in text-independent mode. This form of speaker verification is generally less dependent on time-sequence information, and there is some evidence in a related study of computational efficiency [21] that speaker verification systems might be relatively insensitive to packet loss. One potential anomaly in this hypothesis, equally applicable to both speech and speaker recognitions, is the effect of lost packets on dynamic features which are computed from their static counterparts over some small window, typically in the order of 100 milliseconds or more. Unless appropriately compensated, packet loss of static features would lead to corrupt dynamic features and performance degradation. This difficulty is circumvented here by assuming that the transmitted features are in fact specific to speech and speaker recognitions rather than conventional codec parameters (as defined in the ETSI AURORA standard [22]). As a consequence, packet loss encompasses both static and dynamic features. Preliminary experiments using a Gilbert model (Section 3.1.3) showed very little sensitivity to the patterns of packet loss, so a balanced loss ( $p = 0.25$  and  $q = 0.5$ ) is simulated here with the emphasis placed on the total loss as a percentage of the original.

Experiments are performed with a conventional implementation of a GMM [23] as used by most of today's text-independent speaker verification systems.

##### 4.2. Additive noise

The second degradation considered here typifies the conditions under which wireless mobile devices are commonly used, namely, with a meaningful level of background noise.

The consequences of such additive noise are

- (i) direct contamination of the speech signal,
- (ii) induced changes in the speaking style of the persons subjected to the noise, known as the Lombard reflex [24].

In these experiments, noise is added to the speech recordings thereby minimising any Lombard effects. The noise is added at a moderate level of 15 dB SNR. Subsequently, for completeness, a simple speech enhancement process is applied to the degraded signal.

The form of enhancement considered here has the option of returning the speech to the time domain. Such an approach might lead to suboptimal compensation in terms of recognition performance but nonetheless offers benefits in terms of integration into existing systems and communications networks.

Perhaps the first notable work in this field is that of Boll [25] and Berouti et al. [26] both in 1979. Speech enhancement for human-to-human conversation was performed by an approach still known today as spectral subtraction.

Subsequently, Lockwood and Boudy [27] applied spectral subtraction extensively to automatic speech recognition.

There are many approaches and applications of spectral subtraction. Of particular interest here is an implementation of spectral subtraction termed quantile-based noise

estimation (QBNE), proposed by Stahl et al. [28]. QBNE is an extension of the histogram approach presented by Hirsch and Ehrlicher [29]. The main advantage of these approaches is that an explicit speech, nonspeech detector is not required. Noise estimates are continually updated during both non-speech and speech periods from frequency-dependent, temporal statistics of the degraded speech signal. An efficient implementation of QBNE, important in the context of mobile systems, is described in [30].

### 4.3. Experimental results

#### 4.3.1. Database

The experimental work here was performed on the Speech-Dat Welsh database [20]. The data consists of 2000 speakers recorded over a fixed telephony network. One thousand of the 2000 speakers were used to create a world model and the other 1000 speakers used for speaker model training and testing. Training was performed on approximately 30 seconds of phonetically rich sentences per speaker with a total of about 8 hours for the world model. Two separate text-independent tests used either a 4-digit string, or a single digit, per speaker per test, giving 1000 tests per experiment. Features are standard MFCC-14 static concatenated with 14 dynamic coefficients.

#### 4.3.2. Packet loss and additive noise degradations

To simulate packet loss, approximately 50% of speech features are discarded from the test set, iteratively. No attempt is made to recover these lost vectors although the minimum number of feature vectors per test is capped to two.

Some results are presented in Figures 8 and 9. The detection error trade-off (DET) curves show the system to be highly resilient with minimal increases in error rates until over 75% of the feature vectors are lost, the first three profiles being very close together. This is true for both plots: (Figure 8), the longer, 4-digit string test utterances and (Figure 9) the shorter, single-digit test utterance. Interestingly, in both cases, the profiles diverge toward the left. Considering the 4-digit case (left plot), this indicates that for operating points accepting high false acceptances in return for lower false rejections, the system is particularly robust against packet loss: just 2% false rejections with 50% false acceptances at the extreme case of 98% data loss.

Evidence is presented again in Figure 10 where the EERs are plotted against percentage vector loss and it is clear that the performance begins to degrade only after over 75% of the vectors are lost. This is very much in line with the findings of Section 3 and of McLaughlin et al. [21] who report that a factor of 20 losses can be tolerated before meaningful speaker verification degradation occurs. This finding supports the idea that, in the context of *text-independent* speaker recognition where time sequence information is less critical, there is a large redundancy in typical speech frame rates.

To simulate speaker verification in adverse conditions, the test data is artificially contaminated with car noise at a moderate level of approximately 15 dB SNR.

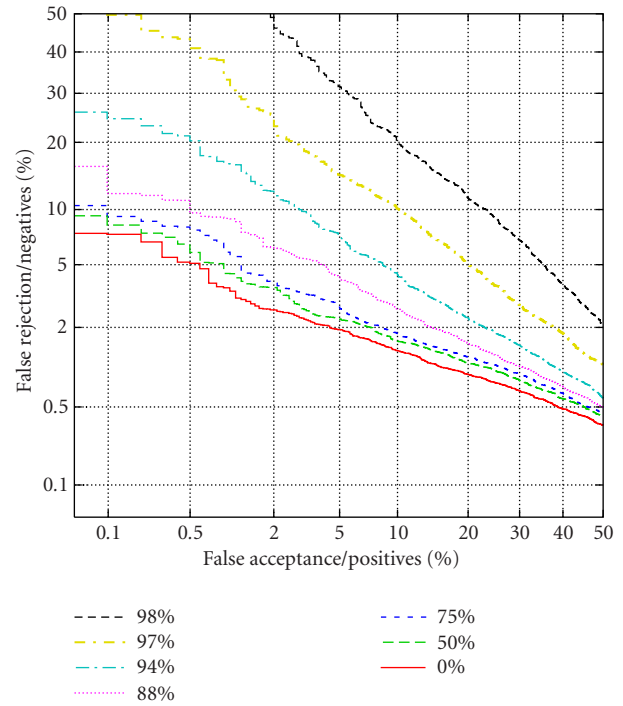


FIGURE 8: Speaker verification performance for varying degrees of feature vector loss, from 0 up to 98% (with a minimum of 2 feature vectors maintained in all tests) for 4-digit string tests.

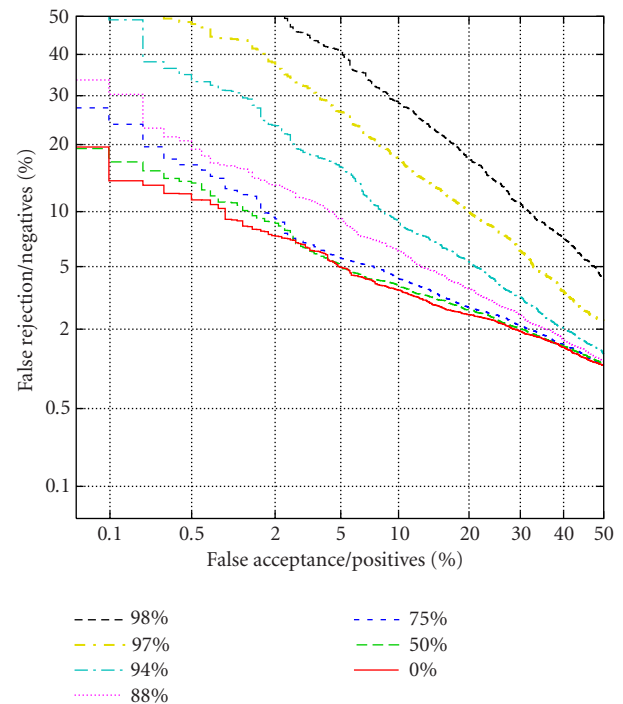


FIGURE 9: Speaker verification performance for varying degrees of feature vector loss, from 0 up to 98% (with a minimum of 2 feature vectors maintained in all tests) for single-digit tests.

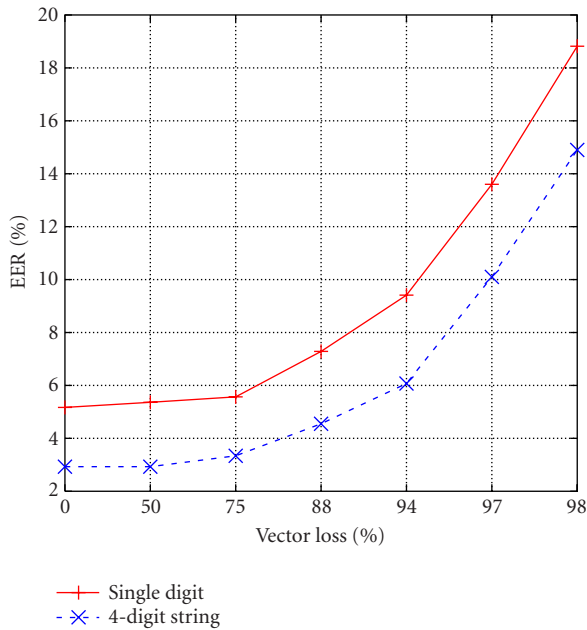


FIGURE 10: EER against feature vector loss (%) for test utterances of 4-digit string (lower profile) and single-digit utterance (upper profile). In all cases, minimum test length is maintained at two vectors.

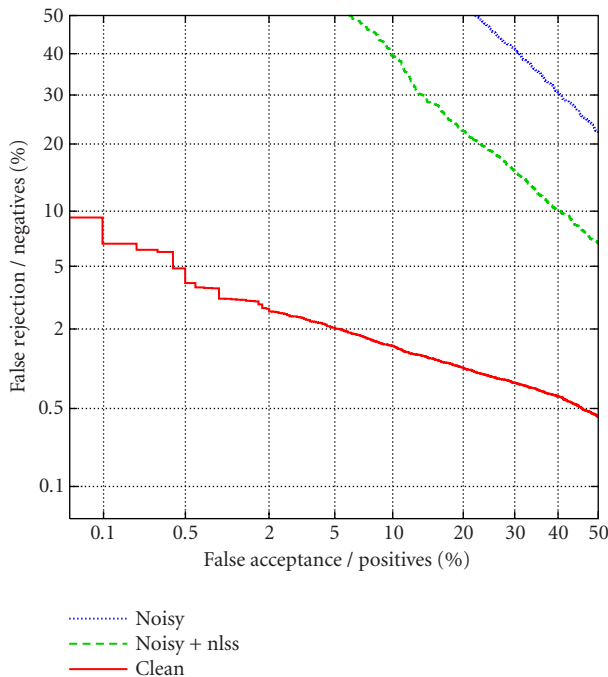


FIGURE 11: Speaker verification performance for the 4-digit string test set with top profile: 15 dB SNR added noise; middle profile: 15 dB SNR added noise plus speech enhancement; and bottom profile: original baseline.

Figure 11 illustrates the effects. The three profiles are for the original telephony test data (bottom profile), the contam-

inated test data (top profile), and the contaminated data after processing with the speech enhancement approach outlined above (middle profile).

Clearly, the levels of performance degradations are marked, even after compensation. This serves to illustrate how relatively small the degradation from packet loss might prove to be in relation to additive noise.

## 5. CONCLUSION

This paper has focused on the emerging need of vocal biometric user authentication over the Internet. More precisely, it has presented the constraints tied with the use of the Internet transmission channel, at the protocol level and at the speech signal level.

At the protocol level, the proposed results have shown that a client-server architecture for vocal biometric user authentication over the Internet involves the transmission of data models or features instead of raw biometric materials. A data encryption process for the client-server link has also been recommended.

At the signal level, the experiments have shown that the packet loss is not a main problem for text-independent vocal person authentication. This is in contrast with previous speech recognition experiments where packet loss was found to reduce the accuracy significantly. Moreover, a large degradation of the performance is observed where a low bit rate coder is used. In this case, packet loss increases the degradation.

Experiments using artificially noised wireless audio records have confirmed that environmental noise remains a main drawback for vocal biometric authentication over the Internet.

## REFERENCES

- [1] J. Abbate, *Inventing the Internet*, Inside Technology. MIT Press, Cambridge, Mass, USA, 2000.
- [2] D. M. Piscitello and A. L. Chapin, "Introduction to routing," *Connexions Magazine*, vol. 7, no. 9, pp. 66–73, 1993.
- [3] J. Siau and A. M. Ariyaeinia, "Biometrics over the internet," COST 275 Technical Meeting, INST Paris, April 2002, <http://www.fub.it/cost275>.
- [4] B. Schneier, "Description of a new variable-length key, 64-bit block cipher (blowfish)," in *Fast Software Encryption, Cambridge Security Workshop, December 1993*, pp. 191–204, Springer-Verlag, Berlin, 1994.
- [5] National Bureau of Standards, "Data Encryption Standard," FIPS Publication 46, U.S. Department of Commerce, Washington, DC, USA, 1977.
- [6] L. Besacier, C. Bergamini, D. Vaufreydaz, and E. Castelli, "The effect of speech and audio compression on speech recognition performance," in *Proc. IEEE Workshop on Multimedia Signal Processing*, Cannes, France, October 2001.
- [7] U. Black, *Voice over IP*, Prentice Hall, Upper Saddle River, NJ, USA, 2001.
- [8] J.-C. Bolot, S. Fosse-Parisis, and D. Towsley, "Adaptive FEC-based error control for internet telephony," in *Proc. 18th Annual Joint Conference of the IEEE Computer and Communications Societies*, pp. 1453–1460, New York, NY, USA, March 1999.

- [9] M. Yajnik, S. Moon, J. Kurose, and D. Towsley, "Measurement and modelling of the temporal dependence in packet loss," in *Proc. 18th Annual Joint Conference of the IEEE Computer and Communications Societies*, pp. 345–352, New York, NY, USA, March 1999.
- [10] The ELISA consortium, "The ELISA systems for the NIST'99 evaluation in speaker detection and tracking," *Digital Signal Processing*, vol. 10, no. 1–3, pp. 143–153, 2000.
- [11] I. Magrin-Chagnolleau, G. Gravier, and R. Blouet for the ELISA consortium, "Overview of the 2000–2001 ELISA consortium research activities," in *2001: A Speaker Odyssey—The Speaker Recognition Workshop*, pp. 67–72, Crete, Greece, June 2001.
- [12] D. A. Reynolds, "Speaker identification and verification using gaussian mixture speaker models," *Speech Communication*, vol. 17, no. 1–2, pp. 91–108, 1995.
- [13] J.-L. Gauvain and C.-H. Lee, "Maximum a posteriori estimation for multivariate Gaussian mixture observations of Markov chains," *IEEE Trans. Speech, and Audio Processing*, vol. 2, no. 2, pp. 291–298, 1994.
- [14] K.-P. Li and J. E. Porter, "Normalizations and selection of speech segments for speaker recognition scoring," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 1, pp. 595–598, New York, NY, USA, April 1988.
- [15] R. Auckenthaler, M. Carey, and H. Lloyd-Thomas, "Score normalization for text-independent speaker verification systems," *Digital Signal Processing*, vol. 10, no. 1, pp. 42–54, 2000.
- [16] M. Ben, R. Blouet, and F. Bimbot, "A Monte-Carlo method for score normalization in automatic speaker verification using Kullback-Leibler distances," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, Orlando, Fla, USA, May 2002.
- [17] P. Mayorga-Ortiz, R. Lamy, and L. Besacier, "Recovering of packet loss for distributed speech recognition," in *Proc. 11th European Signal Processing Conference*, Toulouse, France, September 2002.
- [18] C. Perkins, O. Hodson, and V. Hardman, "A survey of packet loss recovery techniques for streaming audio," *IEEE Network Magazine*, vol. 12, no. 5, pp. 40–48, 1998.
- [19] B. Milner and S. Semnani, "Robust speech recognition over IP networks," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, Istanbul, Turkey, June 2000.
- [20] R. J. Jones, J. S. D. Mason, R. O. Jones, L. Helliker, and M. Pawlewski, "SpeechDat Cymru: A large-scale Welsh telephony database," in *Proc. 1st International Conference on Language Resources and Evaluation: Workshop on Language Resources for European Minority Languages*, Granada, Spain, May 1998.
- [21] J. McLaughlin, D. A. Reynolds, and T. Gleason, "A study of computation speed-ups of the GMM-UBM speaker recognition system," in *Proc. 6th European Conference on Speech Communication and Technology*, vol. 3, pp. 1215–1218, Budapest, Hungary, September 1999.
- [22] D. Pearce, "Enabling new speech driven services for mobile devices: An overview of the ETSI standards activities for distributed speech recognition front-ends," in *Applied Voice Input/Output Society Conference*, San Jose, Calif, USA, May 2000.
- [23] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted Gaussian mixture models," *Digital Signal Processing*, vol. 10, no. 1–3, pp. 19–41, 2000.
- [24] J.-C. Junqua, "The Lombard reflex and its role on human listeners and automatic speech recognizers," *Journal of the Acoustical Society of America*, vol. 93, no. 1, pp. 510–524, 1993.
- [25] S. F. Boll, "Suppression of acoustic noise in speech using spectral subtraction," *IEEE Trans. Acoustics, Speech, and Signal Processing*, vol. 27, no. 2, pp. 113–120, 1979.
- [26] M. Berouti, R. Schwartz, and J. Makhoul, "Enhancement of speech corrupted by acoustic noise," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, pp. 208–211, Washington, DC, USA, April 1979.
- [27] P. Lockwood and J. Boudy, "Experiments with a non-linear spectral subtractor (NSS), hidden Markov models and the projection, for robust speech recognition in cars," in *Proc. 2nd European Conference on Speech Communication and Technology*, vol. 1, pp. 79–82, Genova, Italy, September 1991.
- [28] V. Stahl, A. Fischer, and R. Bippus, "Quantile based noise estimation for spectral subtraction and wiener filtering," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 3, pp. 1875–1878, Istanbul, Turkey, June 2000.
- [29] H. G. Hirsch and C. Ehrlicher, "Noise estimation techniques for robust speech recognition," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 1, pp. 153–156, Detroit, Mich, USA, May 1995.
- [30] N. W. D. Evans, J. S. Mason, and B. Fauve, "Efficient real-time noise estimation without explicit speech, non-speech detection: an assessment on the AURORA corpus," in *Proc. 14th International Conference on Digital Signal Processing*, vol. 2, pp. 985–988, Swansea, Wales, UK, July 2002.

**Laurent Besacier** received his Ph.D. degree in computer science in April 1998 on "A parallel model for automatic speaker recognition" from the University of Avignon, France. Then he spent one and a half year at IMT (Switzerland) as an Associate Researcher working on M2VTS European project (multimodal person authentication). Since September 1999, he has been an Associate Professor in the University of Joseph Fourier (Grenoble). He carries out research on automatic speech and speaker recognition within the GEOD team at CLIPS Lab. He published about 30 papers on various aspects of speech recognition and speaker recognition. He is in the board of AFCEP, the French Speaking Speech Communication Association. He is the supervisor of over 5 Ph.D. students in the area of speaker and speech recognition. His research interests lie in automatic speech and speaker recognition: indexation and tracking of audio documents.



**Aladdin M. Ariyaeinia** received his B.Eng degree in telecommunication engineering, M.S. in digital signal processing, and Ph.D. degree in artificial intelligence in 1976, 1982, and 1986, respectively. He was awarded Chartered Engineer status in 1988. In 1986, he joined the University of East Anglia as a Senior Research Fellow. Over the last fourteen years, Ariyaeinia has been working in the Faculty of Engineering and Information Sciences, the University of Hertfordshire. During this period, he has been conducting research on various aspects of speech processing in close collaboration with industry. He is now a Reader in signal processing, and responsible for leading the Multimedia and Internet Technologies Group. Ariyaeinia's current research interests include speaker and language recognition, speaker-based audio-visual data indexation, biometrics-based recognition over the Internet, and speech enhancement. He has over 40 publications, and has served on various scientific committees.





**John S. Mason** is a Senior Lecturer at the Department of Electrical and Electronic Engineering. He received his M.S. and Ph.D. degrees from the University of Surrey in 1971 and 1974, respectively, joining the University of Wales Swansea as a Lecturer in May 1973. In 1979, he took up a one-year appointment as a Senior Research Engineer at Hewlett Packard Ltd in South Queensferry, and in 1994 he was invited to work on an international project in the Australian National University, Canberra, as a Visiting Research Fellow. From the time of his Ph.D. studies through today, his research interest has focused on digital signal processing. Of a particular note is the work done on finding solutions to complex Chebyshev approximations, widely acknowledged as the first to solve this long-standing problem. More recently, his research has revolved around speech and speaker recognition and multimedia signal processing. In these areas, he has served on the technical committees of a number of international research meetings.



**Jean-François Bonastre** has been an Associate Professor at the LIA, the University of Avignon computer laboratory since 1994. He studied computer science in the University of Marseille and obtained a DEA (Master) in artificial intelligence in 1990. He obtained his Ph.D. degree in 1994, from the University of Avignon, and his HDR (Ph.D. supervision diploma) in 2000, both in computer science, both on speech science, more precisely, on speaker recognition. J.-F. Bonastre is the current President of the AFCP, the French Speaking Speech Communication Association (a Regional Branch of ISCA). He was the Chairman of the RLA2C workshop (1998) and a member of the Program Committee of Speaker Odyssey Workshops (2001 and 2004). J.-F. Bonastre has been an Invited Professor at Panasonic Speech Technology Lab. (PSTL), Calif, USA, in 2002.



**Pedro Mayorga** received his diploma in physics from Universidad Autonoma de Baja California (UABC), Ensenada, Baja California, in 1992. He received the M.S. degree in digital systems from Instituto Politecnico Nacional de Mexico in 1998 with a thesis on fingerprints recognition using neural networks. He is currently a Ph.D. student at CLIPS laboratory, Institut National Polytechnique de Grenoble, France. His Ph.D. research involves the application of digital signal processing and speech recognition to the vocal recognition servers. From 1988 to 1992, he worked in the computer laboratory, Facultad de Ciencias, UABC, Ensenada, Baja California. From 1993 to 1994, he was an Associate Research and Associate Professor at the Instituto de Ingenieria, UABC, Mexicali, Baja California, at the Department of Electrical and Electronic. He was a teacher in microprocessors and digital control systems courses at the same university department. From 1993 to 1994, he worked as a part-time Associate Professor at the Instituto Tecnológico de Mexicali (ITM).



Since 1994, he has been a whole-time Associate Professor at the Department of Electrical-Electronic, ITM. From 1994 to 2000, he was teaching in microprocessors, microcontrollers, digital control, signal processing, and signal and systems courses at the ITM Engineering School.

**Corinne Fredouille** obtained her Ph.D. degree in 2000 in the field of automatic speaker recognition. She has joined the computer science laboratory LIA, University of Avignon, and more precisely the speech processing team, as an Assistant Professor in 2003. Currently, she is an active member of the European ELISA Consortium, of AFCP, the French Speaking Speech Communication Association, and of ISCA/SIG SPLC (Speaker and Language Characterization Special Interest Group).



**Johann Siau** received his B.Eng degree in electrical and electronic engineering in 1997. In 2001, he joined the University of Hertfordshire as a full-time academic staff. Over the past few years, Johann has been working with the Faculty of Engineering and Information Sciences and during this period, he has been conducting research on various aspects of speaker verification technologies, including verification over the Internet. He is currently managing a significant part of the multimedia network at the department and is a member of the Multimedia and Internet Technologies Group. Johann's current research interests include speaker recognition, biometrics-based recognition over the Internet, and network security and vulnerabilities.



**Nicholas W. D. Evans** received the M.Eng degree in electronics and computing science from the University of Wales Swansea in 1999. He then joined the Speech and Image Research Group to pursue his Ph.D. degree sponsored by the Engineering and Physical Sciences Research Council. In 2002, he became a Lecturer in communications at the School of Engineering. Nick's research interests include time-frequency analysis for noise estimation, speech enhancement, noise compensation, noise robust automatic speech recognition, and biometric speaker verification. Nick is a member of ISCA and IEE.



**Roland Auckenthaler** worked with Enigma Ltd. as a Teaching Company Associate from 1998 to 2000 and received his Ph.D. degree from the University of Wales Swansea in 2002 in the area of speaker verification. He now works with Ubiquity Software Corporation in the area of Internet telephony and does part-time research within the University of Wales Swansea. Roland is also a holder of a patent in the area of speaker verification.





**Robert Stapert** moved from the Netherlands to the UK in 1996. There, in his capacity as a Ph.D. student, he spent four years at Swansea University's Speech and Image Processing laboratory. His theme was enhancing speaker verification by means of time sequence information. He completed his Ph.D. in 2000. Since then, he has been employed at Aculab in Milton Keynes, UK, as a member of their digital signal processing team, working as a Software Engineer. He is responsible for the design and development of Aculab's speaker verification product. Further, he is working on projects related to text to speech, speech recognition, as well as various nonspeech related projects.

