# Telephone Handset Identification by Feature Selection and Sparse Representations

Yannis Panagakis and Constantine Kotropoulos

*Department of Informatics*
*Aristotle University of Thessaloniki*
*Thessaloniki, Greece*
{panagakis,costas}@aiia.csd.auth.gr

*Abstract*—Speech signals convey information not only for the speakers' identity and the spoken language, but also for the acquisition devices used during their recording. Therefore, it is reasonable to perform acquisition device identification by analyzing the recorded speech signal. To this end, the *random spectral features* (RSFs) and the *labeled spectral features* (LSFs) are proposed as intrinsic fingerprints suitable for device identification. The RSFs and the LSFs are extracted by applying unsupervised and supervised feature selection to the mean spectrogram of each speech signal, respectively. State-of-the-art identification accuracy of $97.58\%$ has been obtained by employing LSFs on a set of $8$ telephone handsets, from Lincoln-Labs Handset Database (LLHDB).

## I. Introduction

Speech is the most natural way to communicate between humans. Nowadays, speech communication systems acquire transmit, store, and process the information in digital form. However, the digital speech content can be imperceptibly altered by malicious, even amateur, users who may employ a variety of low-cost audio editing software. This creates a serious threat to the *knowledge life cycle*. Indeed, when hearing is no longer believing, the process of going from data to information, knowledge, understanding and, decision making is severely compromised [1]. The consequences of this threat permeate a wide variety of fields, such as intellectual property, intelligence gathering, forensics, and news reporting to name a few. Currently, theories and tools to combat this threat in the field of *digital speech forensics* are still in their infancy. Moreover, there is an urgent need to advance the state-of-the-art in this field [2].

A first step to remedy the aforementioned threat is to extract forensic evidence about the mechanism involved in the generation of the speech recording by analyzing only the speech signal [2]. That is, to identify the acquisition device by assuming that the acquisition devices along with their associated signal processing chain leave behind *intrinsic traces* in the speech signal. Indeed, the electronic devices, especially when include a microphone, cannot have exactly the same frequency response due to tolerances in the production of the electronic components and the different designs employed by

the various manufacturers [3]. This implies that the recorded speech can be considered as a signal whose spectrum is the product of the genuine speech spectrum, driving the acquisition device, and the frequency response of the latter. Consequently, the recorded speech signal can be exploited in device identification, following a blind-passive approach as opposed to active embedding of watermarks or having access to input-output pairs [2].

Although there are significant advances in image forensics [1], audio forensics are less developed [4]. Few exceptions include the authentication of MP3 [5] and the authentication of speakers' environment [6], [7], [8]. Similarly, a few automatic acquisition device identification systems have been developed. For instance, a method for the classification of $4$ microphones has been proposed in [7]. The speech signal is parameterized by employing time domain features and the mel-frequency cepstral coefficients (MFCCs). The identification of the microphones is performed by a Naive Bayes classifier at a short-time frame level. Accuracies on the order of 60-75% have been reported. In [2], the identification of $8$ landline telephone handsets and $8$ microphones is addressed. In particular, the intrinsic characteristics of the device are captured by a template constructed by appending together the means of a Gaussian mixture trained on the speech recordings of each device. To this end, linear- and mel-scaled cepstral coefficients were employed for speech signal representation. Classification accuracies higher than $90\%$ have been achieved, when a support vector machine (SVM) classifier was employed. Recently, a robust system for the identification of cell-phones has been proposed in [3]. In particular, when the MFCCs extracted from device speech recordings are classified by an SVM, $14$ different cell-phones are identified with an accuracy of $96.42\%$.

In this paper, a novel blind-passive method for landline telephone handset identification is proposed. The method resorts on suitable feature extraction from speech recordings along with feature selection and their sparse representation, enabling to trace the recording device. In particular, two feature selection procedures, one unsupervised and another one supervised, are proposed in order to obtain intrinsic features for tracing the recording device. Given a speech recording, its spectrogram is computed and it is averaged along the time axis, yielding the mean spectrogram. For unsupervised

feature selection, the dimensionality of the mean spectrogram is reduced by random projections [9] yielding the *random spectral features* (RSFs) of speech recording. In the supervised setting, the label information (i.e., the class where each device belongs to) of the training speech recordings is taken into account in order to derive a mapping between the feature space where the mean spectrograms lie onto and the label space. Let the training acquisition devices belong to $K$ different device classes with labels $\mathcal{K} = \{1, 2, \ldots, K\}$. Clearly, the label space has as many dimensions as the device classes are. The mapping between the aforementioned two spaces is obtained by solving a regression problem. This mapping can also be exploited in order to map the test mean spectrogram onto the $K$-dimensional space, which is dominated by the label information. These features are referred to as *labeled spectral features* (LSFs).

The RSFs and LSFs can be used to form overcomplete dictionaries of basis signals for devices' intrinsic traces, which is exploited for *sparse representation-based classification* (SRC) [10]. If sufficient training speech recordings are available for each device, it is possible to express any vector of RSFs or LSFs extracted from an unknown (test) device as a compact linear combination of the dictionary atoms for the device actually used for its recording. This representation is designed to be sparse, because it involves only a small fraction of the dictionary atoms and can be computed efficiently via $\ell_1$-norm optimization. The classification is performed by assigning each vector of test RSFs or LSFs the device identity (ID) the dictionary atoms weighted by non-zero coefficients are associated with.

The performance of the proposed method in the identification of 8 telephone handsets is assessed by conducting experiments on the Lincoln-Labs Handset Database (LLHDB) [11] when a stratified 2-fold cross-validation is applied. For comparison purposes, the mean 23-dimensional MFCC vector of each speech recording is considered as a baseline feature for device characterization. Performance comparisons are made against the linear SVM [12] and the nearest-neighbor (NN) classifier, which employs the cosine similarity measure.

The experimental results demonstrate the effectiveness of the proposed RSFs and LSFs over the MFCCs as device fingerprints, no matter which classifier is employed. Meanwhile, the LSFs are able to achieve an accuracy of 97.58% in device identification, outperforming the state-of-the-art method [2] on the LLHDB dataset.

The paper is organized as follows. In Section 2, the RSFs and the LSFs are introduced and the calculation of the MFCCs is described. The sparse representation-based device identification is detailed in Section 3. The dataset and the experimental results are presented in Section 4. Conclusions are drawn in Section 5.

## II. INTRINSIC DEVICE FINGERPRINT EXTRACTION BY SPECTRAL AND CEPSTRAL FEATURES

The majority of features employed in tasks, such as speech and speaker recognition, spoken language identification, etc.

are based on the spectrum of the speech signal. Assuming that the acquisition device is a linear time-invariant system, the impact of the acquisition device on the recorded speech can be modeled by the convolution of the original speech and the impulse response of the device. Thus, the identity of each acquisition device is embedded into the recorded speech, since the spectrum of any windowed recorded speech segment is the product of the spectrum of the original speech signal and the device frequency response.

Motivated by the aforementioned assumption, the RSFs and the LSFs are proposed here as intrinsic traces of recording devices. These features are derived by applying unsupervised and supervised feature selection to the mean spectrograms of the recordings, respectively. The spectrogram of each recorded speech signal is calculated by employing frames of duration 64 ms with a hop size of 32 ms and 2048 FFT bins. Then, the logarithm of the spectrogram is calculated and averaged along the time axis, yielding a 2048-dimensional mean spectrogram.

The RSFs are obtained as follows. The dimensionality of the mean spectrogram is reduced to $d < 2048$ by employing a $d \times 2048$ orthogonal random Gaussian matrix, as described in [9]. Clearly, random projections can be interpreted as an unsupervised feature selection method, since a number $d$ out of 2048 features is selected for acquisition device representation. Let $\mathbf{X} \in \mathbb{R}^{d \times N}$ be the data matrix that contains $N$ vectors of RSFs of size $d$ in its columns. The entries of $\mathbf{X}$ are further post-processed as follows: Each row of $\mathbf{X}$ is normalized to the range $[0, 1]$ by subtracting from each matrix element the row minimum and then by dividing it with the difference between the row maximum and the row minimum.

In order to extract discriminant features from the mean spectrograms the label information of the devices that belong to the training set is taken into account. In particular, we aim to derive features that are highly dependent on the labels. Clearly, the label space is spanned by the columns of the matrix $\mathbf{L} \in \{0, 1\}^{K \times N}$, where the $k$th component of the $n$th column of $\mathbf{L}$, $\mathbf{l}_n$, is 1 if the $n$th device belongs to class $k \in \mathcal{K}$. Let $\mathbf{X}_t \in \mathbb{R}^{2048 \times N_t}$ be the training data matrix, containing in its columns the 2048-dimensional mean spectrograms extracted from, $N_t$ speech signals recorded by using acquisition devices from $K$ classes. Let also $\mathbf{L}_t \in \{0, 1\}^{K \times N_t}$ be the submatrix of $\mathbf{L}$ that is limited to the training samples. Features highly dependent on the labels can be obtained by seeking a linear mapping $\mathbf{M} \in \mathbb{R}^{K \times 2048}$ such that the space of the training mean spectrograms is mapped onto the label space, i.e., $\mathbf{L}_t = \mathbf{M}\,\mathbf{X}_t$. The aforementioned problem can be casted as a regression problem, since it involves the identification of the relationship between sets of dependent variables and independent ones. Although, a simple least squares regression could be employed to derive $\mathbf{M}$, it is well known that such an approach suffers from overfitting. To remedy this drawback of the least squares regression, $\mathbf{M}$ is found by solving the following ridge regression problem:

$$\underset{\mathbf{M}}{\arg\min} \|\mathbf{L}_t - \mathbf{M}\,\mathbf{X}_t\|_F^2 + \lambda\|\mathbf{M}\|_F^2, \qquad (1)$$

where $\lambda$ is a regularization parameter (e.g., the value $\lambda = 0.5$ was used in the experiments) and $\|.\|_F$ denotes the Frobenius norm. The unique closed form solution of (1) is

$$\mathbf{M} = \mathbf{L}_t \mathbf{X}_t^T \left( \mathbf{X}_t \mathbf{X}_t^T + \lambda \mathbf{I} \right)^{-1}. \tag{2}$$

$\mathbf{I}$ denotes the identity matrix of compatible dimensions. In the test phase, by premultiplying any mean spectrogram by $\mathbf{M}$, the $K$-dimensional vector of the LSFs is obtained.

The MFCCs are considered as baseline features [2]. They encode the frequency content of the speech signal by parameterizing the rough shape of spectral envelope. The success of the MFCCs in device identification is justified in [3]. Roughly speaking, the logarithm which involved in the calculation of the MFCCs is a nonlinear transformation with additive property in the spectrum magnitude domain and thus the cepstral features can be consider as a superposition of latent variables, which are related to the recording device, and variables which, are related to the speech content. Following [2], the MFCC calculation employs frames of duration 20 ms with a hop size of 10 ms, and a 42-band filter bank. The correlation between the frequency bands is reduced by applying the discrete cosine transform along the log-energies of the bands. The sequence of 23-dimensional MFCCs is averaged along the time axis yielding a 23-dimensional mean vector. The data matrix containing the MFCCs is postprocessed as described previously for the RSFs.

In Figs. 1 and 2, the mean spectrograms and the MFCCs are depicted, for the same speech utterance recorded by 8 different telephone handsets, respectively. Clearly, both the mean spectrograms and the MFCCs convey discriminant information for the recording device.

## III. ACQUISITION DEVICE IDENTIFICATION VIA SPARSE REPRESENTATION

The problem of revealing the device identity of a vector of RSFs or LSFs, given a number of labeled RSFs or LSFs, respectively, from $N$ acquisition devices is addressed based on the SRC [10]. In the following, when LSFs are employed, $d = K$.

Let us denote by $\mathbf{A}_i = [\mathbf{a}_{i,1} | \mathbf{a}_{i,2} | \ldots | \mathbf{a}_{i,n_i}] \in \mathbb{R}^{d \times n_i}$ the dictionary that contains $n_i$ either RSFs or LSFs stemming from the $i$th device as column vectors (i.e., dictionary atoms). Given a vector of test RSFs (or LSFs) $\mathbf{y} \in \mathbb{R}^d$ that comes from the $i$th device, we can assume that $\mathbf{y}$ is expressed as a linear combination of the atoms that are associated to the $i$th device, i.e.,

$$\mathbf{y} = \sum_{j=1}^{n_i} \mathbf{a}_{i,j} \, c_{i,j} = \mathbf{A}_i \, \mathbf{c}_i \tag{3}$$

where $c_{i,j} \in \mathbb{R}$ are coefficients, which form the coefficient vector $\mathbf{c}_i = [c_{i,1}, c_{i,2}, \ldots, c_{i,n_i}]^T$.

Next, let $\mathbf{A} = [\mathbf{A}_1 | \mathbf{A}_2 | \ldots | \mathbf{A}_N] \in \mathbb{R}^{d \times n}$ be an overcomplete dictionary formed by concatenating $n$ RSFs (or LSFs), which stem from $N$ acquisition devices. Thus, the linear representation of $\mathbf{y} \in \mathbb{R}^d$ in (3) can be equivalently rewritten as

$$\mathbf{y} = \mathbf{A} \, \mathbf{c} \tag{4}$$

where $\mathbf{c} = [\mathbf{0}^T | \ldots | \mathbf{0}^T | \mathbf{c}_i^T | \mathbf{0}^T | \ldots | \mathbf{0}^T]^T$ is the $n \times 1$ augmented coefficient vector, whose elements are zero except those associated with the $i$th device. Thus, the entries of $\mathbf{c}$ contain information about the device the test vector of RSFs (or LSFs) $\mathbf{y} \in \mathbb{R}^d$ comes from.

Since the device ID of a test vector of RSFs (or LSFs) is unknown, we can predict it by seeking the sparsest solution to the linear system of equations $\mathbf{y} = \mathbf{A} \, \mathbf{c}$. Formally, given the overcomplete dictionary $\mathbf{A}$ and the vector of test RSFs or (LSFs) $\mathbf{y} \in \mathbb{R}^d$, the problem of sparse representation is to find the coefficient vector $\mathbf{c}$, such that $\mathbf{y} = \mathbf{A} \, \mathbf{c}$ and $\|\mathbf{c}\|_0$ is minimized, i.e.,

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{c}\|_0 \quad \text{subject to} \quad \mathbf{A}\mathbf{c} = \mathbf{y} \tag{5}$$

where $\|.\|_0$ is the $\ell_0$ quasi-norm returning the number of the non-zero entries of a vector. Finding the solution of the optimization problem (5) is NP-hard due to the nature of the underlying combinational optimization. An approximate solution to the problem (5) can be obtained by replacing the $\ell_0$ norm with the $\ell_1$ norm:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{c}\|_1 \quad \text{subject to} \quad \mathbf{A} \, \mathbf{c} = \mathbf{y} \tag{6}$$

where $\|.\|_1$ denotes the $\ell_1$ norm of a vector. In [13], it has been proved that if the solution is sparse enough, then the solution of (5) is equivalent to the solution of (6), which can be obtained by standard linear programming methods in polynomial time.

A test vector of RSFs (or LSFs) can be classified as follows. The coefficient vector $\mathbf{c}^*$ is obtained by solving (6). Ideally, $\mathbf{c}^*$ contains non-zero entries in positions associated with the dictionary atoms (i.e., columns of $\mathbf{A}$) stemming from a single device, so that we can easily assign the vector of test RSFs (or LSFs) $\mathbf{y}$ to that device. However, due to modeling errors, there are small non-zero entries in $\mathbf{c}^*$ that are associated to multiple devices. To cope with this problem, each RSF (or LSF) is classified to the device class that minimizes the residual $r_i(\mathbf{y}) = \|\mathbf{y} - \mathbf{A} \, \delta_i(\mathbf{c})\|_2$, where $\delta_i(\mathbf{c}) \in \mathbb{R}^n$ is a new vector, whose nonzero entries are associated to the $i$th device only [10]. It is worth mentioning that, the SRC avoids under-fitting, since it employs multiple training samples (instead of the nearest one in the case of the NN) for each class to linearly extrapolate the test sample, but it uses only the smallest necessary number of them to avoid over-fitting. Furthermore, for each test sample, the number of samples needed is automatically determined, since under mild assumptions the $\ell_1$ norm minimization is equivalent to the $\ell_0$ norm minimization [13]. As a result, the SRC can better exploit the actual distributions of the training samples of each class and and therefore it is likely to be more discriminant than other classifiers.

In Fig. 3 (a), the sparse representation coefficients $\mathbf{c}$ for a test vector of RSFs $\mathbf{y}$ extracted from a carbon-button telephone
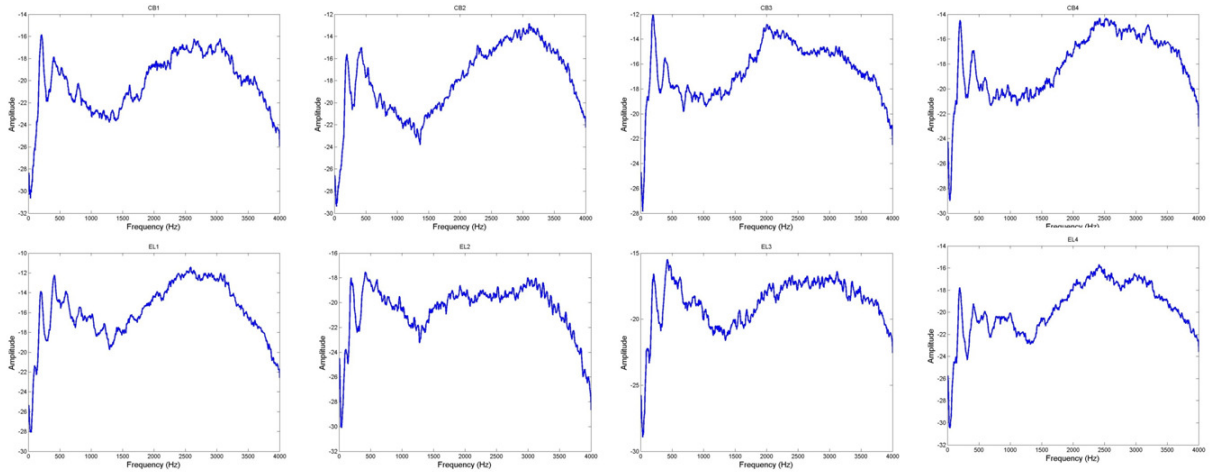
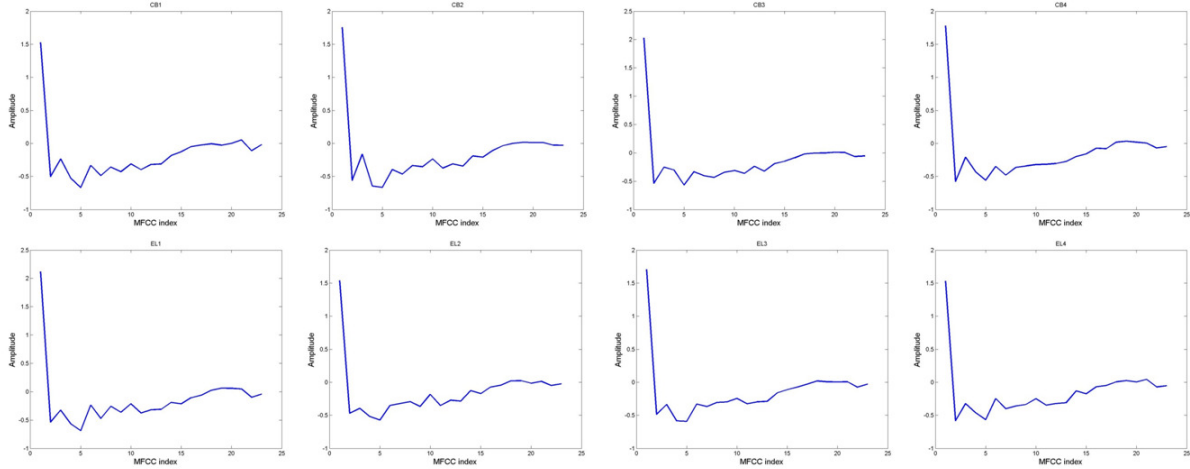Fig. 1.   Mean spectrograms of a speech utterance recorded by 8 different telephone handsets in LLHDB.



Fig. 2.   23-dimensional mean MFCCs of a speech utterance recorded by 8 different telephone handsets in LLHDB.
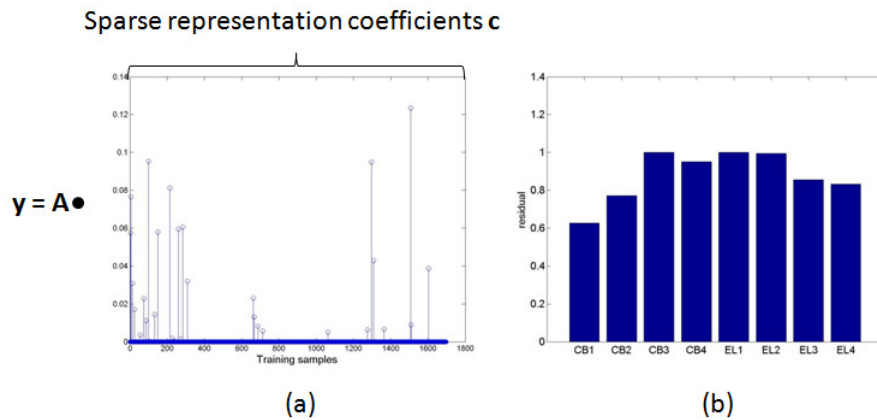


Fig. 3.   The test vector of RSFs $\mathbf{y}$ has been extracted by a carbon-button telephone handset with the ID: CB1. (a) The values of the sparse coefficients $\mathbf{c}$. The non-zero entries of $\mathbf{c}$ are mainly associated with RSFs extracted from speech utterances recorded with the CB1. (b) The residuals $r_i(\mathbf{y})$ of the RSFs. The smallest residual value reveals the identity of the telephone handset (i.e., CB1).

handset with the ID CB1 are illustrated. Fig. 3 (b) shows the     residual $r_i(\mathbf{y})$ with respect to 8 telephone handset IDs.

TABLE I
BEST TELEPHONE HANDSET IDENTIFICATION ACCURACIES ACHIEVED BY THE RSFs, THE LSFs, AND THE MFCCs, WHEN THE SRC, THE LINEAR SVM, AND THE NN ARE EMPLOYED.

| Features | Feature dimension | Classifier | Accuracy (%) |
|---|---|---|---|
| RSFs | 325 | SRC | 95.55 |
| RSFs | 625 | SVM | 94.81 |
| RSFs | 475 | NN | 88.23 |
| LSFs | 8 | SRC | **97.14** |
| LSFs | 8 | SVM | **97.58** |
| LSFs | 8 | NN | 96.52 |
| MFCCs | 23 | SRC | 89.79 |
| MFCCs | 23 | SVM | 87.35 |
| MFCCs | 23 | NN | 81.95 |
| MFCCs- based Gaussian supervector [2] | N/A | SVM | 93.20 |

## IV. EXPERIMENTAL EVALUATION

In order to assess the performance of the proposed method in acquisition device identification, experiments were conducted on the same subset of the Lincoln-Labs Handset Database (LLHDB) [11] as in [2]. This subset consists of speech recordings from 53 speakers (24 males and 29 females) acquired by 8 landline telephone handsets. The first 4 telephone handsets are are carbon-button (CB1-CB4) and the remaining 4 are electrect (EL1-EL4). Following the experimental set-up used in [2], stratified 2-fold cross-validation is employed in the experiments conducted on the LLHDB.

The best identification accuracies are summarized in Table I, when the RSFs, the LSFs, and the MFCCs are classified by the SRC [10], the linear SVM [12], and the NN with the cosine similarity measure. By inspecting Table I, it is clear that the RSFs and the LSFs are able to identify the acquisition device committing less errors than the MFFCs, no matter which classifier is employed. Moreover, the LSFs achieve state-of-the-art identification accuracy if they are fed to either the SVM, or the NN, or the SRC classifier. The SVM achieves the best reported identification accuracy (i.e., $97.58\%$) on the LLHDB.

The performance of the RSFs in telephone handset identification as a function of features dimension (i.e., $d$) is depicted in Fig. 4. It is clear that for $d > 200$ the SRC outperforms the state-of-the-art reported in [2], demonstrating the robustness of the proposed approach in acquisition device identification. The accurate telephone handset identification by the RSFs and their sparse representations is attributed to the following fact. It is well known that by projecting the data onto an orthogonal random Gaussian matrix, the dictionary $\mathbf{A}$ obeys the restricted isometry property (RIP) of a certain, appropriate order (say $S$) [14]. When this property holds, $\mathbf{A}$ approximately preserves the Euclidean length of $S$-sparse RSFs, which in turn implies that $S$-sparse vectors cannot be in the null space of $\mathbf{A}$. The latter is needed, since otherwise there would be no hope for reconstructing these vectors. Clearly, it cannot be guaranteed that the RIP holds for the dictionary constructed by employing the MFCCs as atoms.

The LSFs outperforms both the RSFs and the MFCCs, since they are obtained following a supervised feature selection process. Insight to the performance achieved by LSFs when the SRC, the SVM, and the NN classifier is employed is offered by the confusion matrices shown in Fig. 5, Fig. 6, and Fig. 7, respectively. The rows of the confusion matrices correspond to the predicted device and the columns indicate the actual device. The gray shading in these Figures highlights the fact that most of the identification errors remain within the transducer class (i.e., carbon-button and electrect). The carbon-button telephone handsets are identified more accurately than the electrect ones. This result is attributed to the fact that the transfer functions between the various carbon-button telephone handsets are quite different. Similar results were reported in [2].
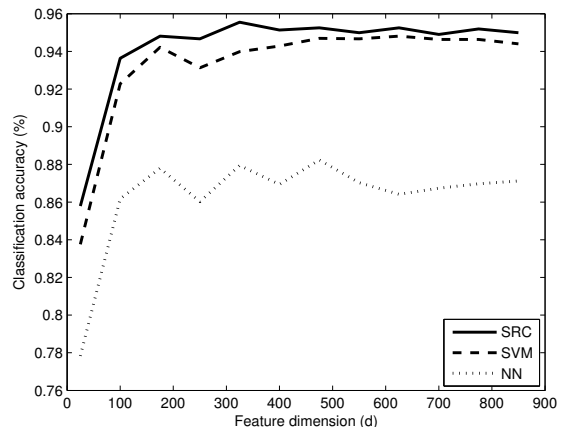


Fig. 4. Telephone handsets identification accuracy for the RSFs obtained by the SRC, the SVM, and the NN on the LLHDB.

## V. CONCLUSIONS

A promising method for telephone handset identification from speech signals has been proposed. The RSFs and the LSFs have been demonstrated to capture the intrinsic trace of the acquisition device, while the sparse representation-based classification has been shown to be able to identify the acquisition device. The experimental results validate the robustness of the RSFs and the LSFs over the MFCCs for
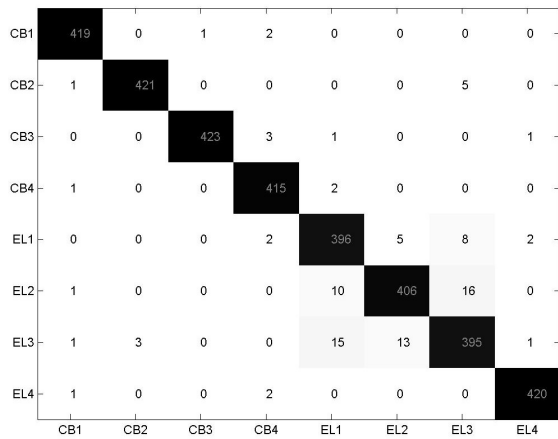
|      | CB1 | CB2 | CB3 | CB4 | EL1 | EL2 | EL3 | EL4 |
|------|-----|-----|-----|-----|-----|-----|-----|-----|
| CB1  | 419 | 0   | 1   | 2   | 0   | 0   | 0   | 0   |
| CB2  | 1   | 421 | 0   | 0   | 0   | 0   | 5   | 0   |
| CB3  | 0   | 0   | 423 | 3   | 1   | 0   | 0   | 1   |
| CB4  | 1   | 0   | 0   | 415 | 2   | 0   | 0   | 0   |
| EL1  | 0   | 0   | 0   | 2   | 396 | 5   | 8   | 2   |
| EL2  | 1   | 0   | 0   | 0   | 10  | 406 | 16  | 0   |
| EL3  | 1   | 3   | 0   | 0   | 15  | 13  | 395 | 1   |
| EL4  | 1   | 0   | 0   | 2   | 0   | 0   | 0   | 420 |

Fig. 5. Confusion matrix for 8 telephone handsets based on the LSFs, when they are classified by the SRC.

|      | CB1 | CB2 | CB3 | CB4 | EL1 | EL2 | EL3 | EL4 |
|------|-----|-----|-----|-----|-----|-----|-----|-----|
| CB1  | 420 | 0   | 1   | 2   | 0   | 0   | 0   | 0   |
| CB2  | 1   | 420 | 0   | 0   | 0   | 1   | 4   | 0   |
| CB3  | 0   | 0   | 421 | 3   | 1   | 0   | 0   | 1   |
| CB4  | 1   | 0   | 1   | 414 | 2   | 0   | 0   | 1   |
| EL1  | 1   | 0   | 0   | 3   | 392 | 6   | 9   | 2   |
| EL2  | 0   | 0   | 0   | 0   | 12  | 399 | 22  | 0   |
| EL3  | 1   | 4   | 0   | 0   | 17  | 18  | 389 | 1   |
| EL4  | 0   | 0   | 1   | 2   | 0   | 0   | 0   | 419 |

Fig. 7. Confusion matrix for 8 telephone handsets based on the LSFs, when they are classified by a NN.

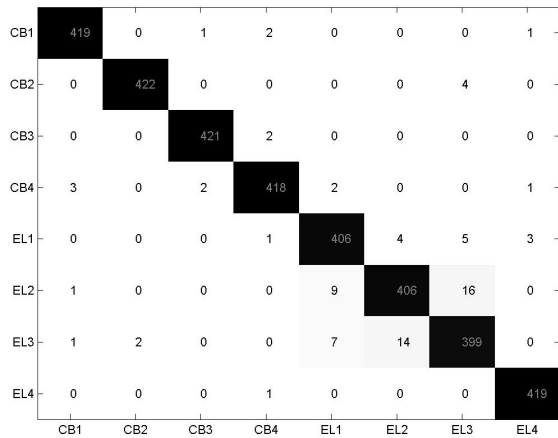|      | CB1 | CB2 | CB3 | CB4 | EL1 | EL2 | EL3 | EL4 |
|------|-----|-----|-----|-----|-----|-----|-----|-----|
| CB1  | 419 | 0   | 1   | 2   | 0   | 0   | 0   | 1   |
| CB2  | 0   | 422 | 0   | 0   | 0   | 0   | 4   | 0   |
| CB3  | 0   | 0   | 421 | 2   | 0   | 0   | 0   | 0   |
| CB4  | 3   | 0   | 2   | 418 | 2   | 0   | 0   | 1   |
| EL1  | 0   | 0   | 0   | 1   | 406 | 4   | 5   | 3   |
| EL2  | 1   | 0   | 0   | 0   | 9   | 406 | 16  | 0   |
| EL3  | 1   | 2   | 0   | 0   | 7   | 14  | 399 | 0   |
| EL4  | 0   | 0   | 0   | 1   | 0   | 0   | 0   | 419 |

Fig. 6. Confusion matrix for 8 telephone handsets based on the LSFs, when they are classified by the SVM.

device characterization, yielding a state-of-the-art performance in recognizing 8 telephone handsets from the LLHDB.

### REFERENCES

[1] H. Farid, "Digital image forensics," *Scientific American*, vol. 6, no. 298, pp. 66–71, 2008.

[2] D. Garcia-Romero and C. Y. Espy-Wilson, "Automatic acquisition device identification from speech recordings," in *Proc. 2010 IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, Dallas, Texas, USA, 2010, pp. 1806–1809.

[3] C. Hanilci, F. Ertas, T. Ertas, and O. Eskidere, "Recognition of brand and models of cell-phones from recorded speech signals," *IEEE Trans. Information Forensics and Security*, vol. 7, no. 2, pp. 625–634, 2012.

[4] R. Maher, "Audio forensic examination," *IEEE Signal Processing Magazine*, vol. 26, no. 2, pp. 84–94, 2009.

[5] R. Yang, Z. Qu, and J. Huang, "Detecting digital audio forgeries by checking frame offsets," in *Proc. 10th ACM Workshop on Multimedia and Security*, New York, NY, USA, 2008, pp. 21–26.

[6] A. Oermann, A. Lang, and J. Dittmann, "Verifier-tuple for audio-forensic to determine speaker environment," in *Proc. 7th ACM Workshop on Multimedia and Security*, New York, NY, USA, 2005, pp. 57–62.

[7] C. Kraetzer, A. Oermann, J. Dittmann, and A. Lang, "Digital audio forensics: a first practical evaluation on microphone and environment classification," in *Proc. 9th ACM Workshop Multimedia and Security*, Dallas, Texas, USA, 2007, pp. 63–74.

[8] H. Malik and H. Farid, "Audio forensics from acoustic reverberation," in *Proc. 2010 IEEE Int. Conf. Acoustics Speech and Signal Processing*, Dallas, Texas, USA, 2010, pp. 1710–1713.

[9] E. Bingham and H. Mannila, "Random projection in dimensionality reduction: applications to image and text data," in *Proc. 7th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Francisco, California, USA, 2001, pp. 245–250.

[10] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 2, pp. 210–227, 2009.

[11] D. Reynolds, "HTIMIT and LLHDB: speech corpora for the study of handset transducer effects," in *Proc. 1997 IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, vol. 2, Munich, Germany, 1997, pp. 1535–1538.

[12] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, pp. 1–27, 2011.

[13] D. Donoho, "For most large underdetermined systems of equations, the minimal l1-norm near-solution approximates the sparsest near-solution," *Communications on Pure and Applied Mathematics*, vol. 59, no. 7, pp. 907–934, 2006.

[14] E. Candes and T. Tao, "Decoding by linear programming," *IEEE Trans. Inform. Theory*, vol. 51, no. 12, pp. 4203–4215, 2005.