

Identification of Sparse Audio Tampering Using Distributed Source Coding and Compressive Sensing Techniques

G. Valenzise, G. Prandi, M. Tagliasacchi, A. Sarti

Abstract

The increasing development of peer-to-peer networks for delivering and sharing multimedia files poses the problem of how to protect these contents from unauthorized and, possibly, malicious manipulations. In the past few years, a large amount of techniques, including multimedia hashes and digital watermarking, have been proposed to identify whether a multimedia content has been illegally tampered or not. Nevertheless, very few efforts have been devoted to identifying which kind of attack has been carried out, with the aim of assessing whether the modified content is still meaningful for the final user and, hopefully, of recovering the original content semantics. One of the main issues that have prevented multimedia hashes from being used for tampering identification is the large amount of data required for this task. Generally the size of the hash should be kept as small as possible to reduce the bandwidth overhead. To overcome this limitation, we propose a novel hashing scheme which exploits the paradigms of compressive sensing and distributed source coding to generate a compact hash signature, and apply it to the case of audio content protection. The audio content provider produces a small hash signature by computing a limited number of random projections of a perceptual, time-frequency representation of the original audio stream; the audio hash is given by the syndrome bits of an LDPC code applied to the projections. At the content user side, the hash is decoded using distributed source coding tools, provided that the distortion introduced by tampering is not too high. If the tampering is sparsifiable or compressible in some orthonormal basis or redundant dictionary (e.g. DCT or wavelet), it is possible to identify the time-frequency position of the attack, with a hash size as small as 200 bits/second: the bit saving obtained by introducing distributed source coding ranges between 20% to 70%.

Part of this work has been presented in the 11th International Conference on Digital Audio Effects, Espoo, Finland, September 2008 [1]. This work has been partially sponsored by the EU under Visnet II Network of Excellence.

The authors are with Dipartimento di Elettronica e Informazione, Politecnico di Milano, P.zza Leonardo da Vinci, 32 20133 - Milano, Italy - Ph. +39-02-2399-7624 - FAX: +39-02-2399-7321 - E-mail: valenzise@elet.polimi.it, prandi@elet.polimi.it, tagliasa@elet.polimi.it, sarti@elet.polimi.it

I. INTRODUCTION

With the increasing diffusion of digital multimedia contents in the last years, the possibility of tampering with multimedia contents – an ability traditionally reserved, in the case of analog signals, to few people due to the prohibitive cost of the professional equipment – has become quite a widespread practice. In addition to the ease of such manipulations, the problem of the diffusion of unauthorized copies of multimedia contents is exacerbated by security vulnerabilities and peer-to-peer sharing over the Internet, where digital contents are typically distributed and posted. This is particularly true for the case of audio files, which represent the most common example of digitally distributed multimedia contents. Some versions of the same audio piece may differ from the original because of processing, due for example to compression, resampling, or transcoding at intermediate nodes. In other cases, however, malicious attacks may occur by tampering with part of the audio stream and possibly affecting its semantic content. Examples of this second kind of attacks are the alteration of a piece of evidence in a criminal trial, or the manipulation of public opinion through the use of false wiretapping. Often, for the sake of information integrity, not only it is useful to detect whether the audio content has been modified or not, but also to identify which kind of attack has been carried out. The reasons why it is generally preferred to identify how the content has been tampered with are twofold: on one hand, given an estimate of *where* the signal was manipulated, one can establish whether or not the audio file is still meaningful for the final user; on the other hand, in some circumstances, it may be possible to recover the original semantics of the audio file.

In the past literature, the aim of distinguishing legitimately modified copies from manipulations of a multimedia file has been addressed with two kinds of approaches: watermarks and media hashes. Both approaches have been extensively applied to the case of image content types, whether fewer systems have been proposed for the case of audio signals. Digital watermarking techniques embed information directly into the media data to ensure both data integrity and authentication. Even if digital watermarks can be categorized based on several properties, such as robustness, security, complexity, invertibility, etc. [2], a common taxonomy is to distinguish between *robust* and *fragile* watermarks. It is the latter category that is particularly useful for checking the integrity of an audio file: a fragile watermark is a mark that is easily altered or destroyed when the host data is modified through some transformation, either legitimate or not. If the watermark is designed to be robust with respect to legitimate, perceptually-irrelevant modifications (e.g. compression or resampling), and at the same time to be fragile with respect to perceptually and semantic significant alterations, then it is a *content-fragile* watermark [2]. With this scheme, a possible

tampering can be detected and localized by identifying the damage to the extracted watermark. Examples of this approach for the case of image content types are given in [3] and [4]. The authors of [5] propose an image authentication scheme that is able to localize tampering, by embedding a watermark in the wavelet coefficients of an image. If a tampering occurs, the system provides information on specific frequencies and space regions of the image that have been modified. This allows the user to make application-dependent decisions concerning whether an image, which is JPEG compressed for instance, still has credibility. A similar idea, also working on the signal wavelet domain, has been applied to audio in [6], with the aim of copyright verification and tampering identification. The image watermarking system devised in [7] inserts a fragile watermark in the least significant bits of the image on a block-based fashion; when a portion of the image is tampered with, only the watermark in the corresponding blocks is destroyed, and the manipulation can be localized. Celik et al. [8] extend this method by inserting the watermark in a hierarchical way, to improve robustness against vector quantization attacks. In [9], image protection and tampering localization is achieved through a technique called “cocktail watermarking”: two complementary watermarks are embedded in the original image to improve the robustness of the detector response, while at the same time enabling tampering localization. The same ideas have been applied by the authors to the case of sounds [10], by inserting the watermark in the host audio FFT coefficients. For a more exhaustive review of audio watermarking for authentication and tampering identification see Steinebach and Dittmann [2].

Despite their widespread diffusion as a tool for multimedia protection, watermarking schemes suffer from a series of disadvantages: 1) watermarking authentication is not backward compatible with previously encoded contents (unmarked contents cannot be authenticated later by just retrieving the corresponding hash); 2) the original content is distorted by the watermark; 3) the bit-rate required to compress a multimedia content might increase due to the embedded watermark. An alternative solution for authentication and tampering identification is the use of *multimedia hashes*. Unlike watermarks, content hashing embeds a signature of the original content as part of the header information, or can provide a hash separately from the content upon a user’s request. Multimedia hashes are inspired by cryptographic digital signatures, but instead of being sensitive to single bit changes, they are supposed to offer proof of perceptual integrity. Despite some audio hashing systems (also named *audio fingerprinting*) have been proposed in the past few years [11][12][13], most of the previous research, as for the case of watermarking, has concentrated on the case of images [14], [15]. In [11], the authors build audio fingerprints by collecting and quantizing a number of robust and informative features from an audio file, with the purpose of audio identification as well as fast database lookup. Haitsma and Kalker [12] build audio fingerprints robust to legitimate

content modifications (mp3 compression, resampling, moderate time and pitch scaling), by dividing the audio signal in highly overlapping frames of about 0.3 seconds: for each frame, they compute a frequency representation of the signal through a filter bank with logarithmic spacing among the bands, in order to resemble the Human Auditory System (HAS). The redundancy of musical sounds is exploited by taking the differences between subbands in the same frame, and between the same subbands in adjacent time instants; the resulting vector is quantized with one bit, and similarities between each short fingerprint are computed through the Hamming distance. By concatenating all the fingerprints of each frame, a global hash is obtained, which is used next to efficiently query a song database of previously encoded fingerprints. Though in principle such approach could be used for identifying possible localized tampering in the audio stream, the authors do not explicitly address this problem. An excellent review of algorithms and applications of audio fingerprinting is presented in [13].

To the best of the authors' knowledge, no audio hashing technique has been used up to now with the purpose of detecting and localizing unauthorized audio tampering. One of the main reasons of that is probably the great amount of bits of the audio hashes required for enabling the identification of the tampering, when traditional fingerprinting approaches as the ones described above are employed. In fact, in order to limit the rate overhead, the size of the hash needs to be as small as possible. At the same time, the goal of tampering localization calls for increasing the hash size, in order to capture as much as possible about the original multimedia object. Recently, Lin et al. have proposed a new hashing technique for authentication [15] and tampering localization [16] for images, which produce very short hashes by leveraging distributed source coding theory: in this system, the hash is composed of the Slepian-Wolf encoding bit-stream of a number of quantized random projections of the original image; the content user computes its own random projections on the received (and possibly tampered) image, and uses them as a side information to decode the received hash. By setting some maximum pre-defined tampering level on the received image (e.g. a minimum tolerated PSNR between the original and the forged image is allowed), it is possible to transmit the hash without the need of a feedback channel, performing rate allocation at the encoder side (a similar bit allocation technique has been adopted by the authors also in the context of reduced-reference image quality assessment [17]). When decoding succeeds, it is possible to identify tampered regions of the image, at the cost of additional hash bits. This scheme has been applied also to the case of audio files [18]: instead of random projections of pixels, the authors compute for each signal frame a weighted spectral flatness measure, with randomly chosen weights, and encode this information to obtain the hash. Though this scheme applies well to the authentication task (which can be attained with a hash overhead less than 100 bits/second), it is not clear how to extend the application

to identification of general kinds of tampering.

We have recently proposed a new image hashing technique [19] which exploits both the distributed source coding paradigm and the recent developments in the theory of compressive sensing. The algorithm proposed in this paper extends these ideas to the scenario of audio tampering. It also shares some similarities with the works in [16] and [18]: as in [18], the hash is generated by computing random projections starting from a perceptually-significant time-frequency representation of the audio signal and storing the syndrome bits obtained by LDPC (Low Density Parity-Check Codes) encoding the quantized coefficients. With respect to [18], the proposed algorithm is novel in the following aspect: by leveraging compressive sensing principles, we are able to identify tamperings that are not sparse in the time domain only, but that can be represented by a sparse set of coefficients in some orthonormal basis or redundant dictionary. Even if the spatial models introduced in [16] could be thought of as a representation of the tampering in some dictionary, it is apparent that the compressive sensing interpretation allows much more flexibility in the choice of the sparsifying basis, since it just uses off-the-shelf basis expansions (e.g. wavelet or DCT) which can be added to the system for free.

To clear up which are the capabilities and the limitations of the proposed system, Figure 1 shows an example of malicious tampering with an audio signal. This demonstration has been carried out on a piece of audio speech, with a length of approximately 2 seconds, read from a newspaper by a speaker. The whole recording, which is about 32 seconds long, has also been used as a proof of concept to present some experimental results on the system in Section VII. Figure 1(a) shows the original waveform, which corresponds to the Italian sentence: “Un sequestro da tredici milioni di euro” (“A confiscation of thirteen million euros”). This sentence has been tampered with in order to substitute the words “tredici milioni” (“thirteen millions”) with “quindici miliardi” (“fifteen billions”), see Figure 1(b). In order to compute the hash, as explained in Section IV, we compute a coarse-scale perceptual time-frequency map of the signal (in this case, with a temporal resolution of 1/4 seconds). From the received tampered waveform and from the information of the hash, the user is able to identify the tampering (Figure 1(d)).

The rest of the paper is organized as follows: Section II provides the necessary background information about compressive sensing and distributed source coding; Section III describes the tampering model; Section IV gives a detailed description of the system; Section VI describes how it is possible to estimate the rate of the hash at the encoder without feedback channel or training; the tampering identification algorithm is tested against various kinds of attacks in Section VII, where also the different bit-rate requirements for the hash using or not distributed source coding are compared; finally, Section VIII draws some concluding remarks.

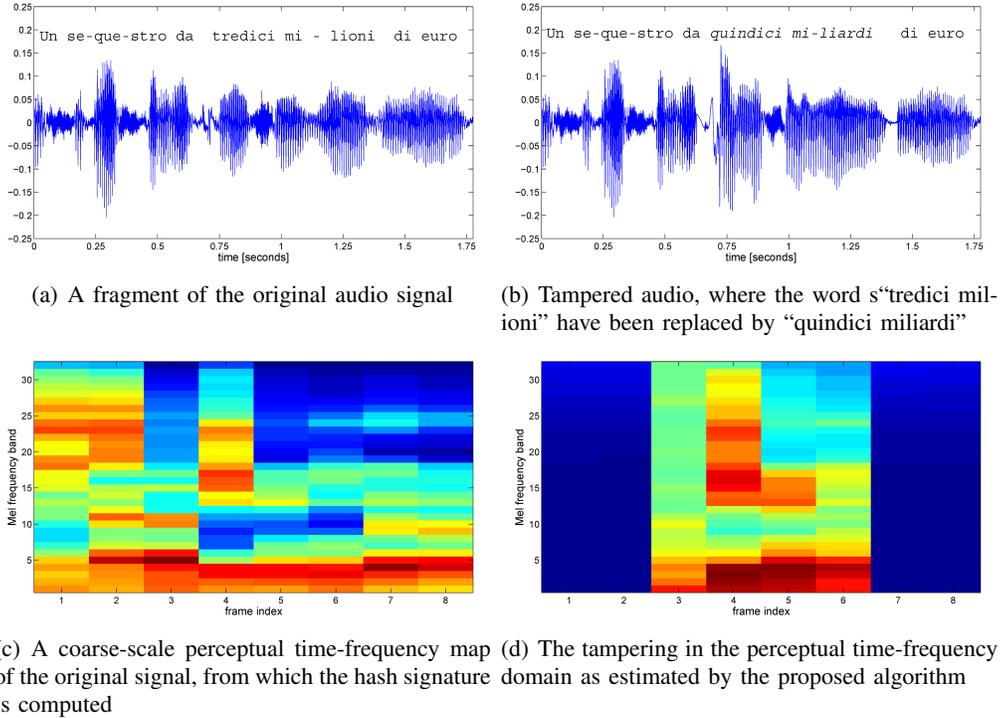


Fig. 1. An example of the result of the proposed audio tampering identification, applied to a fragment of speech read from a newspaper.

II. BACKGROUND

In this section we review the important concepts behind compressive sensing and distributed source coding, that constitute the underlying theory of the proposed tampering identification system. In spite of the relatively large amount of literature published on these fields in the past few years, this is a very concise introduction; for a more detailed and exhaustive explanation the interested reader may refer to [20][21][22] for compressive sensing and to [23][24][25] for distributed source coding.

A. Compressive sampling (CS)

Compressive sampling (or compressed sensing) is a new paradigm which asserts that it is possible to perfectly recover a signal from a limited number of incoherent, non-adaptive linear measurements, provided that the signal admits a sparse representation in some orthonormal basis or redundant dictionary, i.e. it can be represented by a small number of non-zero coefficients in some basis expansion. Let $\mathbf{x} \in \mathbb{R}^n$ be the signal to be acquired, and $\mathbf{y} \in \mathbb{R}^m$, $m < n$, a number of linear random projections (measurements) obtained as $\mathbf{y} = \mathbf{A}\mathbf{x}$. In general, given the prior knowledge that \mathbf{x} is k -sparse, i.e. that only k out of its

n coefficients are different from zero, one can recover \mathbf{x} by solving the following optimization problem:

$$\min \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x} \quad (1)$$

where $\|\cdot\|_0$ simply counts the number of non-zero elements of \mathbf{x} . This program can correctly recover a k -sparse signal from $m = k + 1$ random samples [26]. Unfortunately, such a problem is NP hard, and it is also difficult to solve in practice for problems of moderate size.

To overcome this exhaustive search, the compressive sampling paradigm uses special measurement matrices \mathbf{A} that satisfy the so-called *Restricted Isometry Property* (RIP) of order k [22], which says that all subsets of k columns taken from \mathbf{A} are in fact nearly orthogonal or, equivalently, that linear measurements taken with \mathbf{A} approximatively preserve the Euclidean length of k -sparse signals. This in turn implies that k -sparse vectors cannot be in the null space of \mathbf{A} , a fact that is extremely useful, as otherwise there would be no hope of reconstructing these vectors. Merely verifying that a given \mathbf{A} has the RIP according to the definition is combinatorially complex; however, there are well known cases of matrices that satisfy the RIP, obtained for instance by sampling i.i.d. entries from the normal distribution with mean 0 and variance $1/n$. When the RIP holds, then the following linear program gives an accurate reconstruction:

$$\min \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x}. \quad (2)$$

The solution of (2) is the same as the one of (1) provided that the number of measurements satisfy $m \geq C \cdot k \log_2(n/k)$, where C is some small positive constant. Moreover, if \mathbf{x} is not exactly sparse, but it is at least *compressible* (i.e. its coefficients decay as a power law), then solving (2) guarantees that the quality of the recovered signal is as good as if one knew ahead of time the location of the k largest values of \mathbf{x} and decided to measure those directly [22]. These results also hold when the signal is not sparse as is, but it has a sparse representation in some orthonormal basis. Let $\Psi \in \mathbb{R}^{n \times n}$ denote an orthonormal matrix, whose columns are the basis vectors. Let us assume that we can write $\mathbf{x} = \Psi\alpha$, where α is a k -sparse vector. Clearly, (2) is a special case of this instance, when Ψ is the identity matrix. Given the measurements $\mathbf{y} = \mathbf{A}\mathbf{x}$, the signal \mathbf{x} can be reconstructed by solving the following problem:

$$\min \|\alpha\|_1 \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\Psi\alpha. \quad (3)$$

Problem (3) can be solved without prior knowledge of the actual sparsifying basis Ψ for different test bases, until a sparse reconstruction α is obtained.

In most practical applications, measurements are affected by noise (e.g. quantization noise). Let us

consider noisy measurements $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{z}$, where \mathbf{z} is a norm-bounded noise, i.e. $\|\mathbf{z}\|_2 \leq \epsilon$. An approximation of the original signal \mathbf{x} can be obtained by solving the modified problem:

$$\min \|\boldsymbol{\alpha}\|_1 \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\boldsymbol{\Psi}\boldsymbol{\alpha}\|_2 \leq \epsilon. \quad (4)$$

Problem (4) is an instance of a second order cone program (SOCP) [27] and can be solved in $O(n^3)$ time. Several fast algorithms have been proposed in the literature that attempt to find a solution to (4). In this work, we adopt the SPGL1 algorithm [28], which is specifically designed for large scale sparse reconstruction problems.

B. Distributed source coding (DSC)

Consider the problem of communicating a continuous random variable X . Let Y denote another continuous random variable correlated to X . In a distributed source coding setting, the problem is to decode X to its quantized reconstruction \hat{X} given a constraint on the distortion measure $D = E[d(X, \hat{X})]$ when the side information Y is available only at the decoder. Let us denote by $R_{X|Y}(D)$ the rate-distortion function for the case when Y is also available at the encoder, and by $R_{X|Y}^{WZ}(D)$ the case when only the decoder has access to Y . The Wyner-Ziv theorem [24] states that, in general, $R_{X|Y}^{WZ}(D) \geq R_{X|Y}(D)$ but $R_{X|Y}^{WZ}(D) = R_{X|Y}(D)$ for Gaussian memoryless sources and mean square error (MSE) as distortion measure.

The Wyner-Ziv theorem has been applied especially in the area of video coding under the name of distributed video coding (DVC), where the source X (pixel values or DCT coefficients) is quantized with 2^J levels, and the J bitplanes are independently encoded, computing parity bits by means of a turbo encoder. At the decoder, parity bits are used together with the side information Y to “correct” Y into a quantized version \hat{X} of X , performing turbo decoding, typically starting from the most significant bitplanes. To this end, the decoder needs to know the joint p.d.f. (probability density function) $p_{XY}(X, Y)$. More recently, LDPC codes have been adopted instead of turbo codes [29][30].

Although the rate-distortion performance of a practical DSC codec strongly depends on the actual implementation employed, it is yet possible to approximately quantify the gain obtained by introducing a Wyner-Ziv coding paradigm, in order to estimate the bit saving produced in the hash signature. Let X and Y be zero-mean, i.i.d. Gaussian variables with variance, respectively, σ_X^2 and σ_Y^2 ; also, let σ_N^2 be the variance of the innovation noise $N = Y - X$. Classical information theory [31] asserts that the rate expressed in bits per sample for a given distortion level D , in the case of a Gaussian source X is given

by:

$$R_X(D) = \frac{1}{2} \log_2 \frac{\sigma_X^2}{D}. \quad (5)$$

The rate-distortion function for the case of Wyner-Ziv encoding, when the conditions of the theorem are satisfied, is

$$R_{X|Y}^{WZ}(D) = \frac{1}{2} \log_2 \frac{\sigma_X^2 \sigma_N^2}{D(\sigma_X^2 + \sigma_N^2)} \quad (6)$$

which becomes, in the hypothesis that $\sigma_X^2 \gg \sigma_N^2$, approximately equal to the rate needed to encode the innovation N :

$$R_{X|Y}^{WZ}(D) \approx \frac{1}{2} \log_2 \frac{\sigma_N^2}{D}. \quad (7)$$

Subtracting (7) from (5), we obtain the expected coding gain due to Wyner-Ziv coding:

$$\Delta R_{WZ} = \frac{1}{2} \log_2 \frac{\sigma_X^2}{\sigma_N^2} \quad (8)$$

As we shall see in Section IV, σ_X^2 relates to the energy of the original signal, while σ_N^2 to the energy of the tampering. Equation (8) shows that the advantage of using a DSC approach with respect to a traditional quantization and encoding becomes consistent when the signal and the side information are well correlated, i.e. when the energy of the tampering is small relative to the energy of the original sound.

III. TAMPERING MODEL

Before describing in more detail the architecture of the system, we need to set up a model for sparse tampering. Let $\mathbf{x} \in \mathbb{R}^n$ be the original signal; we model the effect of a sparse tampering $\mathbf{e} \in \mathbb{R}^n$ as

$$\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{e}, \quad (9)$$

where $\tilde{\mathbf{x}}$ is the modified signal received by the user. We postulate without loss of generality that \mathbf{e} has only k non-zero components (in fact, it suffices for \mathbf{e} to be sparse or compressible in some basis or frame).

Let $\mathbf{y} = \mathbf{A}\mathbf{x}$ be the random measurements of the original signal and $\tilde{\mathbf{y}} = \mathbf{A}\tilde{\mathbf{x}}$ be the projections of the tampered signal: clearly, the relation between the tampering and the measurements is given by

$$\mathbf{b} = \tilde{\mathbf{y}} - \mathbf{y} = \mathbf{A}(\tilde{\mathbf{x}} - \mathbf{x}) = \mathbf{A}\mathbf{e}. \quad (10)$$

If the sensing matrix \mathbf{A} is chosen such that it satisfies the RIP, we have that:

$$\|\mathbf{b}\|_2 = \|\mathbf{A}\mathbf{e}\|_2 \approx \sqrt{\frac{m}{n}}\|\mathbf{e}\|_2 \quad (11)$$

and thus we are able to approximate the energy of the tampering from the projections computed at the decoder and the encoder-side projections reconstructed exploiting the hash. This fact comes out to be very useful to estimate the energy of the tampering at the content user side and will be exploited in Section IV. Furthermore in order to apply the Wyner-Ziv theorem, we need \mathbf{b} to be i.i.d. Gaussian with zero mean. This has been verified through experimental simulations on several tampering examples. Indeed, a theoretical justification can be provided by invoking the central limit theorem, since each element $b_i = \sum_{j=1}^n A_{ij}e_j$ is the sum of random variables whose statistics are not explicitly modeled.

IV. DESCRIPTION OF THE SYSTEM

The proposed tampering detection and localization scheme is depicted in Figure 2. The general architecture of the system is composed by two actors: on one hand, there is the *content producer* (CP), which is the entity that publishes or distributes the legitimate and authentic copies of the original audio content. On the other hand, there is the *content user* (CU), which is the consumer of the audio content released by the CP. The CP disseminates copies of the original content $\mathbf{X} \in \mathbb{R}^N$, where N is the total number of audio samples of the signal, through possibly untrusted intermediaries, which may tamper with the authentic file manipulating its semantics; at the same time, the CU may get its own copy $\tilde{\mathbf{X}}$ of the audio file from nodes different from the starting CP. In order to protect the integrity of the multimedia content, the CP builds a small hash signature \mathcal{H} of the audio signal. To perform content authentication, the user sends a request for the hash signature to an authentication server, which is supposed to be trustworthy. By exploiting the hash, the user can estimate the distortion of the received content $\tilde{\mathbf{X}}$ with respect to the original \mathbf{X} . Furthermore, if the tampering is sparse in some basis expansion, the system produces a tampering estimation $\hat{\mathbf{e}}$ which identifies the attack in the time-frequency domain. In the following, we detail the hash generation procedure at the content producer side and the tampering identification at the content user side.

A. Generation of the hash signature

At the content producer side, given the audio stream \mathbf{X} and a random seed S , the encoder generates the hash signature $\mathcal{H}(\mathbf{X}, S)$ as follows:

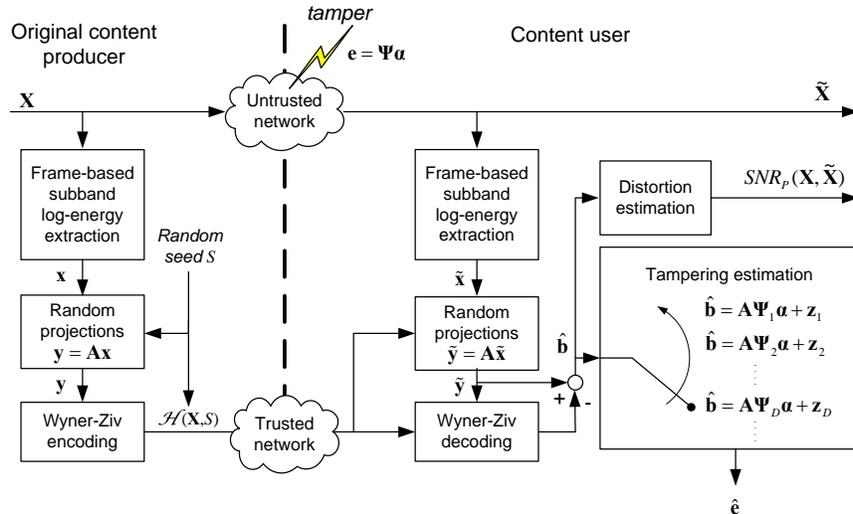


Fig. 2. Block diagram of the proposed tampering identification scheme

1) *Frame based subband log-energy extraction*: The original single-channel audio stream \mathbf{X} is partitioned into non-overlapping frames of length F samples. The power spectrum of each frame is subdivided into U Mel frequency subbands [32], and for each subband the related spectral log-energy is extracted. Let $h_{f,u}$ be the energy value for the u -th band at frame f . The corresponding log-energy value is computed as follows:

$$x_{f,u} = \log(1 + h_{f,u}). \quad (12)$$

The values $x_{f,u}$ provide a time-frequency perceptual map of the audio signal (see Figure 1). The log-energy values are “rasterized” as a vector $\mathbf{x} \in \mathbb{R}^n$, where $n = UN/F$ is the total number of log-energy values extracted from the audio stream.

2) *Random projections*: A number of linear random projections $\mathbf{y} \in \mathbb{R}^m$, $m < n$, is produced as $\mathbf{y} = \mathbf{A}\mathbf{x}$. The entries of the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ are sampled from a Gaussian distribution $\mathcal{N}(0, 1/n)$, using some random seed S , which will be sent as part of the hash to the user.

3) *Wyner-Ziv encoding*: The random projections \mathbf{y} are quantized with a uniform scalar quantizer with step size Δ . As mentioned in Section I, to reduce the number of bits needed to represent the hash, we do not send directly the quantization indices. Instead, we observe that the random projections computed from the possibly tampered audio signal will be available at the decoder side. Therefore, we can perform lossy encoding with side information at the decoder, where the source to be encoded is \mathbf{y} and the “noisy” random projections $\tilde{\mathbf{y}} = \mathbf{A}\tilde{\mathbf{x}}$ play the role of the side information. The vector $\tilde{\mathbf{x}}$ contains the log-

energy values of the audio signal received at the decoder. With respect to the distributed source coding setting illustrated in Section II-B, we have $X = \mathbf{y}$, $Y = \tilde{\mathbf{y}}$, $N = \mathbf{b} = \tilde{\mathbf{y}} - \mathbf{y}$. Following the approach widely adopted in the literature on distributed video coding [25], we perform bitplane extraction on the quantization bin indices. Then each bitplane vector is LDPC coded to create the hash.

B. Hash decoding and tampering identification

The content user receives the (possibly tampered) audio stream $\tilde{\mathbf{X}}$ and requests the syndrome bits and the random seed of the hash $\mathcal{H}(\mathbf{X}, S)$ from the authentication server. On each user's request, a different seed S is used in order to avoid that a malicious attack could exploit the knowledge of the nullspace of \mathbf{A} [15].

1) *Frame-based subband log-energy extraction*: A perceptual, time-frequency representation of the signal $\tilde{\mathbf{X}}$ received by the content user is computed using the same algorithm described above for the content producer side. At this step, the vector $\tilde{\mathbf{x}}$ is produced.

2) *Random projections*: A set of m linear random measurements $\tilde{\mathbf{y}} = \mathbf{A}\tilde{\mathbf{x}}$ are computed using a pseudo-random matrix \mathbf{A} whose entries are drawn from a Gaussian distribution with the same seed S as the encoder.

3) *Wyner-Ziv decoding*: A quantized version $\hat{\mathbf{y}}$ is obtained using the hash syndrome bits and $\tilde{\mathbf{y}}$ as side information. LDPC decoding is performed starting from the most significant bitplane.

- If a feedback channel is available, decoding always succeeds, unless an upper bound is imposed on the maximum number of hash bits.
- Conversely, if the actual distortion between the original and the tampered signal is higher than the maximum tolerated distortion determined by the original content producer, decoding might fail.

4) *Distortion estimation*: If Wyner-Ziv decoding succeeds, an estimate of the distortion in terms of a perceptual signal-to-noise ratio is computed using the projections of the subsampled energy spectrum of the tampering. Let $\hat{\mathbf{b}} = \tilde{\mathbf{y}} - \hat{\mathbf{y}}$ be the projections of the subsampled energy spectrum of the tampering; we define the perceptual signal-to-noise ratio (SNR_P) of the received audio stream as

$$\text{SNR}_P = 10 \log_{10} \frac{\|\hat{\mathbf{y}}\|_2^2}{\|\hat{\mathbf{b}}\|_2^2} \quad [\text{dB}]. \quad (13)$$

This definition needs some further interpretation. In fact, we compute the SNR_P from the projections in place of the whole time-frequency perceptual map of both the signal and the tampering. This is justified by the energy conservation principle stated in (11) and by the fact that, at the content user side, no information about the authentic audio content is available; hence this is an approximation of the actual

SNR_P , which uses the quantized projections obtained by decoding the hash signature, in the reasonable hypothesis that $\|\hat{\mathbf{y}}\| \approx \|\mathbf{y}\|$ and $\|\hat{\mathbf{b}}\| \approx \|\mathbf{b}\|$.

5) *Tampering estimation*: If the tampering can be represented by a sparse set of coefficients in some basis Ψ_i , it can be reconstructed starting from the random projections $\mathbf{b} = \tilde{\mathbf{y}} - \hat{\mathbf{y}}$ by solving the following optimization problem, as anticipated in Section II-A:

$$\min \|\alpha\|_1 \quad \text{s.t.} \quad \|\hat{\mathbf{b}} - \mathbf{A}\Psi_i\alpha\|_2 \leq \epsilon \quad (14)$$

For a given orthonormal basis Ψ_i , the expansion of the tampering in that basis, i.e. $\alpha_i = \Psi_i^T(\mathbf{x} - \tilde{\mathbf{x}})$, might not be sparse enough with respect to the number of available random projections m and the optimization algorithm might not converge to a feasible solution. In such cases, it is not possible to perform tampering identification, and a different orthonormal basis Ψ_j , $j \neq i$ is tested. If the optimization algorithm does not converge for any of the tested bases, the tampering is declared to be non-sparse. This is the case, for example, of quantization noise introduced by audio compression. If the reconstruction succeeds for more than one basis, we choose the one in which the tampering is the sparsest. While, in principle, this just means that we should take the basis that returns the smallest ℓ_0 metrics, we have in practice to cope with reconstruction noise, which in fact prevents the recovered tampering to be exactly sparse. A simple solution is to select the basis that gives the smallest ℓ_1 norm; however, this approach has the drawback of being too sensitive towards high values of the coefficients (e.g. due to different dynamic ranges in the transform domains). As experimentally shown in Section VII-B, this bias has the side-effect that selecting the minimum ℓ_1 norm reconstruction does not ensure that one is performing the best possible tampering estimation. A more effective heuristic is to use some ℓ^p metrics, with $0 < p < 1$, or similar norms, as the ones devised in [33]. In our experiments, we have computed the norm of the coefficients α as

$$\|\alpha\| = \sum_{i=1}^m \arctan\left(\frac{|\alpha_i|}{\delta}\right) \quad (15)$$

where δ has been set so that $\arctan(1/\delta) = 1$.

V. CHOICE OF THE HASH PARAMETERS

In the hash construction procedure, there are two parameters that influence the quality of tampering estimation: the number of random projections m used to build the hash, and the number of bitplanes J which determines the distortion due to quantization on the reconstructed measurements at the user-side. In this section we analyze the trade-off between the rate needed to encode the hash, which also depends on the maximum allowed tampering level as explained in Section VI, and the accuracy of the tampering

estimation: a larger number of bitplanes J and of measurements m correspond to a higher quality of tampering estimation, and at the same time to a higher rate spent for the hash. In order to find an optimal tradeoff between m and J , we conducted Montecarlo simulations on a generic sparse signal \mathbf{x} , with two different sparsity levels k/n . We evaluate the goodness of the tampering estimation by calculating the reconstruction normalized MSE (NMSE_R) between the original k -sparse signal \mathbf{x} and its approximation $\hat{\mathbf{x}}$ obtained by solving problem (4):

$$\text{NMSE}_R = \frac{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2}. \quad (16)$$

The noise $\mathbf{z} = \hat{\mathbf{x}} - \mathbf{x}$ in (4) in this case corresponds to quantization noise, which is uniformly distributed between $-\Delta/2$ and $\Delta/2$, where Δ is the quantization step size. We measure the impact of quantization noise by measuring the signal-to-quantization noise ratio

$$\text{SNR}_y = 10 \log_{10} \frac{\|\mathbf{y}\|_2^2}{\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2}, \quad (17)$$

where $\hat{\mathbf{y}}$ is the quantized version of the random projections $\mathbf{y} = \mathbf{A}\mathbf{x}$. As for the reconstruction basis, Ψ , we just assign $\Psi = I$ in (4), i.e. we assume that the signal is sparse as is, or equivalently that some oracle has told us the optimal sparsifying basis in advance. Figure 3 shows the NMSE_R contour set for two levels of sparsity ($k/n = 0.15$ and $k/n = 0.25$) as a function of the number of projections m and of the quantization distortion of the measurements (SNR_y). We observe a graceful improvement of the performance by increasing either m or SNR_y . For the same values of the parameters, the normalized MSE of the reconstructed signal is lower for sparser signals ($k/n = 0.15$). This is justified by the CS result on the number of projections which requires $m \geq C \cdot k \log_2(n/k)$ (see Section II-A): thus the contour set for $k/n = 0.25$ appears as it was “shifted” to the right with respect to the case $k/n = 0.15$ in Figure 3. As for the quantization of the projections, provided that the number of measurements is compatible with the sparsity level as explained before, we can observe that the value of NMSE_R decreases as SNR_y becomes larger. In a practical scenario, the quantization step size Δ should be chosen in such a way to attain $\text{SNR}_y \geq 25$ dB, in order to be robust with the choice of m , which depends on the actual sparsity of the tampering and on the constant C and is therefore unknown at the content producer side. In our experiments in the rest of the paper, we have set $C = 1.3$.

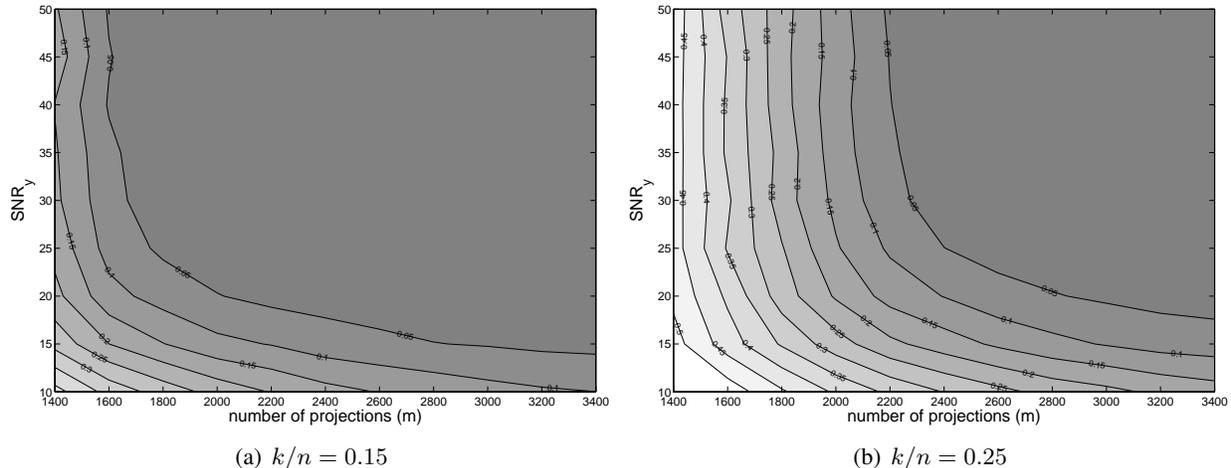


Fig. 3. Normalized MSE of the reconstructed tampering as a function of the number of measurements m and the measures signal-to-quantization noise ratio SNR_y , expressed in dB.

VI. RATE ALLOCATION

In Section III we have shown that the correlation model between the original and the tampered random projections can be written as

$$\tilde{\mathbf{y}} = \mathbf{y} + \mathbf{b} \quad (18)$$

Hereafter we assume that \mathbf{y} and \mathbf{b} are statistically independent. This is reasonable if the tampering is considered independent from the original audio content.

Let $j = 1, \dots, J$ denote the bitplane index and R_j the bitrate (in bits/sample) needed to decode the j -th bitplane. As mentioned in Section III, the probability density function of \mathbf{y} and \mathbf{b} can be well approximated to be zero mean Gaussian, respectively with variance σ_y^2 and σ_b^2 . The rate estimation algorithm receives in input the source variance σ_y^2 , the correlation noise variance σ_b^2 , the quantization step size Δ and the number of bitplanes to be encoded J and returns the average number of bits needed to decode each bitplane R^j , $j = 1, \dots, J$. The value of σ_y^2 can be immediately estimated from the random projections at the time of hash generation. The value of σ_b^2 is set to be equal to the maximum MSE distortion between the original and the tampered signal, for which tampering identification can be attempted.

The rate allocated to each bitplane is given by:

$$R^j = H(\mathbf{y}^j | \tilde{\mathbf{y}}, \mathbf{y}^{j-1}, \mathbf{y}^{j-2}, \dots, \mathbf{y}^1) \text{ [bits/sample]} + \Delta R \quad (19)$$

where \mathbf{y}^j denotes the j -th bitplane of \mathbf{y} . In fact LDPC decoding of bitplane j exploits the knowledge of

	SNR _P [dB]	Sparsity (k/n)	m/n
T	20.3	9 % (1D-DCT)	0.54
F	11.5	26 % (2D-DCT)	0.66
TF	14.5	6 % (Haar)	0.54

TABLE I
PERCEPTUAL SNR, SPARSITY FACTOR k/n IN THE MOST “SPARSIFYING” BASIS (IN PARENTHESES) AND m/n RATIO FOR THE THREE CONSIDERED TAMPERING EXAMPLE.

the real-valued side information $\tilde{\mathbf{y}}$ as well as previously decoded bitplanes $\mathbf{y}^{j-1}, \mathbf{y}^{j-2}, \dots, \mathbf{y}^1$. Since we use nonideal channel codes with a finite sequence length m to perform source coding a rate overhead of approximately $\Delta R = 0.1$ [bit/sample] is added. The integral needed to compute the value of the conditional entropy in (19) is factored out in detail in our previous work [34].

VII. EXPERIMENTAL RESULTS

We have carried out some experiments on 32 seconds of speech audio data, sampled at 44100 Hz and 16 bits per sample. The test audio consists of a piece of a newspaper article read by a speaker; the recording is clean but for some noise added at a few time instants, including the high frequency noise of a shaken key ring, the wide-band noise of some crumpling paper, and some impulsive noise in the form of coughs of the speaker. We have set the size of the audio frame to $F = 11025$ samples (0.25 seconds), and the number of Mel frequency bands to $U = 32$, obtaining a total of 128 audio frames corresponding to $n = 4096$ log-energy coefficients. We have then assembled a testbed considering 3 kinds of tampering:

- *Time localized tampering (T)*: we have replaced some words in the speech at different positions, for a total tampering length of 3.75 seconds (about 11.7% of the total length of the audio sequence);
- *Frequency localized tampering (F)*: a low-pass phone-band filter (cut-off frequency at 3400 Hz and stop frequency at 4000 Hz) is applied to the entire original audio stream;
- *Time-frequency localized tampering (TF)*: a cough at the beginning of the stream and the noise of the key ring in the middle are canceled out using the standard noise removal tool of the “Audacity” free audio editing software [35]. The noise removal tool implemented in this application is an adaptive filter, whose frequency response depends on the local frequency characteristics of the noise. In this case, the total time length of the attack is 4.36 seconds.

The reconstruction of the tampering has been attempted in 3 different bases, besides the log-energy domain: 1-D DCT (Discrete Cosine Transform across frequency bands of the same frame: this corresponds to extracting Mel Frequency Cepstral Coefficients), 2-D DCT (across time and frequency), and 2-D Haar

wavelet. Table I summarizes the perceptual SNRs and the sparsity of the three tampering examples, in the domain where its values is the lowest. It also reports the number of computed projections m in terms of the ratio m/n . Note that this ratio is always less than one (i.e. $m < n$), thus the adopted setting is coherent with the compressive sensing framework explained in Section II-A. In the following, we evaluate two aspects of the system, namely: 1) the rate spent for Wyner-Ziv encoding the hash with respect to the rate that would have been spent for encoding and transmitting the projections without DSC; 2) the relation between the ℓ_1 and the inverse tangent norms of the quality of the reconstructed tampering in different domains.

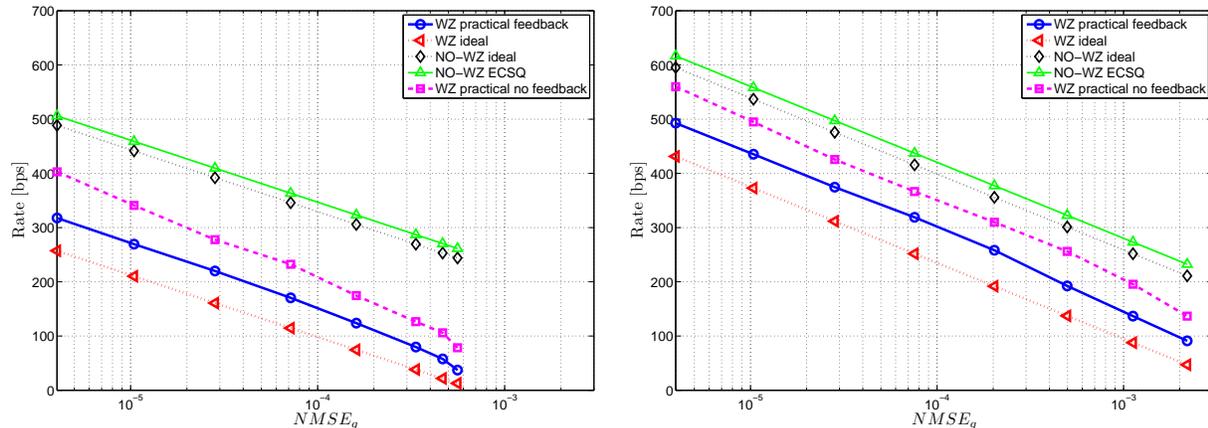
A. Rate-Distortion performance of the hash signature

As described in Section IV, we use distributed source coding for reducing the payload due to the hash. In this section, we want to quantify the bit-saving obtained with Wyner-Ziv coding of the hash. In order to do so, we have compared the rate distortion function of Wyner-Ziv (WZ) coding and of hash direct quantization and transmission, i.e. without using DSC (NO-WZ). Figure 4 depicts these two situations for the cases of the frequency and time domain tampering. In both the two graphs, the value of quantization MSE has been normalized by the energy of the measurements \mathbf{y} , in order to make the result comparable with other possible manipulations:

$$\text{NMSE}_q = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2}{\|\mathbf{y}\|_2^2} \quad (20)$$

The bold dotted lines represents the theoretical WZ rate-distortion curve of the measurements stated in (7). The bold solid and dashed lines represent instead the actual rate-distortion behavior obtained by using a practical WZ codec, either using the feedback channel or directly estimating at the encoder side the rate as explained in Section VI. For comparison, we have also plotted the rate-distortion functions of an ideal NO-WZ uniform quantizer (Shannon's bound), drawn as a thin dotted line, and the rate-distortion curve of an Entropy-Constrained Scalar Quantization (ECSQ), which is a well-studied and effective practical quantization scheme (thin solid line).

We can make two main comments on the curves in the two graphs of Figure 4. The first difference between the frequency and the time tampering is that all the rate-distortion functions in the frequency attack are shifted upwards to higher rates, and have a steeper descending slope as the distortion increases. This is due to the fact that the frequency manipulation has a higher sparsity coefficient k/n , i.e. more measurements are needed for signal reconstruction. Although in the real application no guess about the sparsity of the tampering can be made at the content producer side, here we have fixed a different sparsity



(a) Time sparse tampering, with a sparsity factor k/n set to 0.15. (b) Frequency sparse tampering, with sparsity factor $k/n = 0.25$.

Fig. 4. Rate-distortion function of the hash signature with different encoding approaches.

for the two kinds of attacks, in order to visually prove the effect of the number of measures on the hash length. Thus, even if the rate per measurement is the same in both the cases (it only depends on the signal energy, as expressed in (5) and (7)), the rate in bits per second has slopes and offsets proportional to the number of measurements m . Clearly, if we did not use compressive sensing to reduce the dimensionality of the data (i.e. $\mathbf{y} = \mathbf{x}$ in our setting), the rate required for the hash would have been equivalent to using random projections with $m = n$; therefore, the rate saving due to compressive sensing is approximately equal to the ratio m/n . The second interesting remark that emerges from Figure 4 is the different gap between the family of WZ rates (ideal, with feedback and without feedback) and the NO-WZ curves. As (8) suggests, the coding gain from NO-WZ to WZ strongly depends on the energy of the tampering, i.e. to SNR_P (see Table I). In the case of time attack, we have $\text{SNR}_P^T = 20.3$ dB, while $\text{SNR}_P^F = 11.5$ dB, thus according to (8) the bit saving achieved with WZ is smaller in the case of the frequency attack. As can be inferred from the graphs, this gain ranges from 20% to 70%.

B. Choice of the best tampering reconstruction

In practice, the tampering may be sparse or compressible in more than one basis: this may be the case, for instance, of piece-wise polynomials signals which are generally sparse in several wavelet expansions. When this situation occurs, multiple tampering reconstructions are possible, and at the content user side there is an ambiguity about what is the best tampering estimation. As described in Section IV-B, we are ultimately interested in finding the sparsest tampering representation. This requires in practice to evaluate

	Log-energy	1D-DCT	2D-DCT	Haar Wavelet
T	$7.1 \cdot 10^{-3}$	$4.8 \cdot 10^{-3}$	$2.5 \cdot 10^{-2}$	$7.9 \cdot 10^{-3}$
F	$1.1 \cdot 10^{-1}$	$3.6 \cdot 10^{-2}$	$8.6 \cdot 10^{-3}$	$1.3 \cdot 10^{-2}$
TF	$2.3 \cdot 10^{-3}$	$1.5 \cdot 10^{-3}$	$4.3 \cdot 10^{-3}$	$1.4 \cdot 10^{-3}$

TABLE II
NMSE_R FOR TAMPERING RECONSTRUCTION WITH A HASH AT A BIT RATE OF 200 BPS.

the sparsity of the tampering in each basis expansion: we use for this purpose the inverse-tangent norm defined in (15). To validate the choice of this norm, we compare the optimal basis expansion predicted from the ℓ_1 norm and the inverse tangent norm with the actual best basis in terms of ℓ_2 reconstruction quality.

We evaluate the goodness of the tampering estimation by calculating the reconstruction normalized MSE between the log-energy spectrum of the original tampering and the log-energy spectrum of the estimated one:

$$\text{NMSE}_R = \frac{\|\hat{\mathbf{e}} - \mathbf{e}\|_2^2}{\|\mathbf{e}\|_2^2}. \quad (21)$$

Reconstruction NMSE values obtained with a fixed bit rate for the hash are shown in Tables II (for 200 bps) and III (for 400 bps). The bit rate depends on the number of measurements m (given in Table I) and on the number of bitplanes per measurement J . For a resulting rate of 200 bps, the number of bitplanes for the three kinds of attack (T, F, TF) is, respectively, 7, 5 and 6. When the rate is 400 bps, we have $J = 10$ for the time attack, $J = 8$ for the frequency attack, and $J = 9$ for the time-frequency tampering. From the tables it is clear that, by looking for a sparse tampering in other bases besides the canonical one (log-energy), better results can be achieved using the same hash length, as highlighted by the bold numbers in the tables. In particular, it can be observed that the wide-band, time-localized tampering is better reconstructed using the 1D-DCT basis, which is able to capture tampering correlations only along the frequency axis, avoiding tampering discontinuities over time. The frequency-localized tampering is better reconstructed using the 2D-DCT basis, due to its time extension and wide-band characterization which exhibits only a single discontinuity along the frequency axis. Finally, Haar wavelet is a good compromise to detect time-frequency localized tampering because it is able to deal with discontinuities along both time and frequency axes.

Tables IV and V show the ℓ_1 norms of the reconstructed tampering coefficients in the four analyzed bases. Note that at a rate equal to 200 bps, the ℓ_1 norm suggests, for the time-frequency (TF) tampering, that the best reconstruction is with the 1D-DCT coefficients. However, Table II indicates that the best

	Log-energy	1D-DCT	2D-DCT	Haar Wavelet
T	$2.4 \cdot 10^{-4}$	$2.1 \cdot 10^{-4}$	$1.6 \cdot 10^{-2}$	$4.5 \cdot 10^{-4}$
F	$9.3 \cdot 10^{-2}$	$1.2 \cdot 10^{-2}$	$1.9 \cdot 10^{-3}$	$3.1 \cdot 10^{-3}$
TF	$4.7 \cdot 10^{-5}$	$6.0 \cdot 10^{-5}$	$1.1 \cdot 10^{-3}$	$4.5 \cdot 10^{-5}$

TABLE III
NMSE_R FOR TAMPERING RECONSTRUCTION WITH A HASH AT A BIT RATE OF 400 BPS.

	Log-energy	1D-DCT	2D-DCT	Haar Wavelet
T	265.33	183.06	366.05	248.26
F	1219.12	1005.34	251.41	488.08
TF	509.88	256.00	445.95	260.71

TABLE IV
 ℓ_1 -NORM OF THE TAMPERING USING A FIXED BIT RATE FOR THE HASH SIGNATURE OF 200 BPS.

reconstruction is actually in the Haar wavelet domain. This is due to the noise introduced by compressive sensing recovery at low rates, which makes the use of the ℓ_1 norm as an estimator of the sparsity more error-prone. This effect is partially alleviated using the inverse tangent norm, as shown in Tables VI and VII.

To have a visual insight of the effect of different bases in the tampering reconstruction, we have drawn in Figure 5 the log-energy spectrum of the original audio signal and of the frequency-localized (F) tampering, followed by the log-energy spectrum of the tampering reconstructed in two different domains using a hash rate of 200 bps. It is apparent from the figure that the quality of the estimated tampering reconstructed using 2D-DCT considerably overcomes the one obtained in the log-energy domain.

VIII. CONCLUSIONS

We presented a hash-based tampering identification system for detecting and identifying illegitimate manipulations in audio files. The algorithm works with sparse modifications, leveraging the recent compressive sensing results for reconstructing the tampering from a set of random non-adaptive measurements. Perhaps the most distinctive feature of the proposed system is its ability to reconstruct a tampering that is sparse in some orthonormal basis or frame, without knowing at the content producer side the actual content alteration. In practice, such an approach is feasible only if the bit length of the hash is not too large: we have found that encoding the hash signature through a distributed source coding paradigm enables a consistent reduction of the transmitted bits, especially when the strength of the tampering is small compared to the original signal energy. The hash size may be further decreased in the future

	Log-energy	1D-DCT	2D-DCT	Haar Wavelet
T	344.57	246.97	543.06	338.81
F	1761.76	1394.95	445.42	731.58
TF	594.46	330.28	639.89	325.64

TABLE V
 ℓ_1 -NORM OF THE TAMPERING USING A FIXED BIT RATE FOR THE HASH SIGNATURE OF 400 BPS.

	Log-energy	1D-DCT	2D-DCT	Haar Wavelet
T	270.36	166.68	455.10	252.26
F	1115.59	793.15	187.00	323.35
TF	324.44	150.30	349.18	136.60

TABLE VI
 INVERSE TANGENT NORM OF THE TAMPERING USING A FIXED BIT RATE FOR THE HASH SIGNATURE OF 200 BPS.

by considering weighted ℓ_1 minimization [33] to reduce the number of measurements required by the algorithm.

REFERENCES

- [1] G. Prandi, G. Valenzise, M. Tagliasacchi, and A. Sarti, "Detection and identification of sparse audio tampering using distributed source coding and compressive sensing techniques," in *Proc. 11th Int. Conf. on Digital Audio Effects*, Espoo, Finland, 2008.
- [2] M. Steinebach and J. Dittmann, "Watermarking-Based Digital Audio Data Authentication," *EURASIP Journal on Applied Signal Processing*, vol. 2003, no. 10, pp. 1001–1015, 2003.
- [3] J. Fridrich, "Image watermarking for tamper detection," in *IEEE International Conference on Image Processing*, Chicago, October 1998, vol. 2.
- [4] J.J. Eggers and B. Girod, "Blind watermarking applied to image authentication," in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, Salt Lake City, 2001, vol. 3.
- [5] D. Kundur and D. Hatzinakos, "Digital watermarking for telltale tamper proofing and authentication," *Proceedings of the IEEE*, vol. 87, no. 7, pp. 1167–1180, 1999.
- [6] R. Tu and J. Zhao, "A novel semi-fragile audio watermarking scheme," *Proc. of the 2nd IEEE Internatioal Workshop on Haptic, Audio and Visual Environments and Their Applications*, pp. 89–94, 2003.
- [7] P.W. Wong, "A public key watermark for image verification and authentication," *IEEE International Conference on Image Processing, 1998*, vol. 1, 1998.
- [8] M.U. Celik, G. Sharma, E. Saber, and A.M. Tekalp, "Hierarchical Watermarking for Secure Image Authentication With Localization," *IEEE Trans. Image Process.*, vol. 11, no. 6, pp. 585, 2002.
- [9] C.S. Lu, S.K. Huang, C.J. Sze, and H.Y.M. Liao, "Cocktail watermarking for digital image protection," *IEEE Transactions on Multimedia*, vol. 2, no. 4, pp. 209–224, 2000.

	Log-energy	1D-DCT	2D-DCT	Haar Wavelet
T	324.41	224.92	675.91	334.40
F	1586.11	1087.59	412.81	575.68
TF	308.54	196.59	536.68	171.91

TABLE VII
INVERSE TANGENT NORM OF THE TAMPERING USING A FIXED BIT RATE FOR THE HASH SIGNATURE OF 400 BPS.

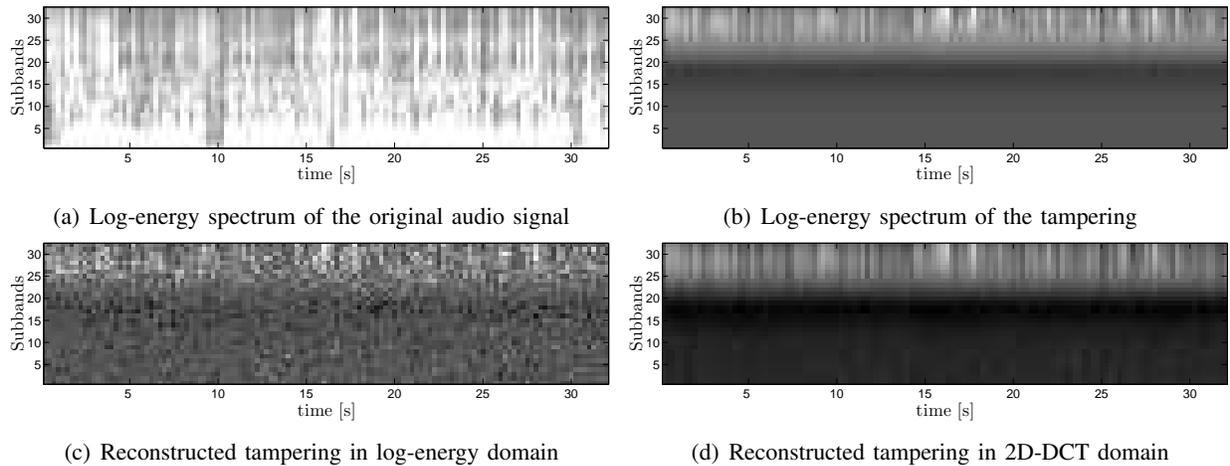


Fig. 5. Example of frequency tampering. The hash length is 200 bps.

- [10] C.S. Lu, H.Y.M. Liao, and L.H. Chen, "Multipurpose audio watermarking," in *Proc. 15th Int. Conf. on Pattern Recognition*, 2000.
- [11] M.K. Mihcak and R. Venkatesan, "A Perceptual Audio Hashing Algorithm: A Tool For Robust Audio Identification and Information Hiding," *Proc. of the 4th Information Hiding Workshop*, vol. 2137, pp. 51–65, 2001.
- [12] J. Haitsma and T. Kalker, "A Highly Robust Audio Fingerprinting System With an Efficient Search Strategy," *Journal of New Music Research*, vol. 32, no. 2, pp. 211–221, 2003.
- [13] P. Cano, E. Batlle, T. Kalker, and J. Haitsma, "A review of algorithms for audio fingerprinting," *Int. Workshop on Multimedia Signal Processing*, 2002.
- [14] S. Roy and Q. Sun, "Robust Hash for Detecting and Localizing Image Tampering," in *IEEE International Conference on Image Processing*, S.Antonio, TX, 2007, vol. 6.
- [15] Y.C. Lin, D. Varodayan, and B. Girod, "Image authentication based on distributed source coding," in *IEEE International Conference on Image Processing*, S.Antonio, TX, September 2007, vol. 3.
- [16] Y.C. Lin, D. Varodayan, and B. Girod, "Spatial Models for Localization of Image Tampering Using Distributed Source Codes," in *Picture Coding Symposium (PCS)*, Lisbon, Portugal, November 2007.
- [17] K. Chono, Y. C. Lin, D. Varodayan, Y. Miyamoto, and B. Girod, "Reduced-reference image quality estimation using distributed source coding," in *IEEE International Conference on Multimedia and Expo*, Hannover, Germany, June 2008.
- [18] D. Varodayan, Y.C. Lin, and B. Girod, "Audio authentication based on distributed source coding," in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, Las Vegas, NV, 2008.

- [19] M. Tagliasacchi, G. Valenzise, and S. Tubaro, "Localization of sparse image tampering via random projections," in *Proc. IEEE Int. Conf. Image Processing*, San Diego, CA, USA, 2008.
- [20] E. Candes, "Compressive sampling," in *International Congress of Mathematicians*, Madrid, Spain, 2006.
- [21] R.G. Baraniuk, "Compressive Sensing," *Signal Processing Magazine, IEEE*, vol. 24, no. 4, pp. 118–121, 2007.
- [22] E.J. Candes and M.B. Wakin, "An introduction to compressive sampling: A sensing/sampling paradigm that goes against the common knowledge in data acquisition," *Signal Processing Magazine, IEEE*, vol. 25, no. 2, pp. 21–30, March 2008.
- [23] D. Slepian and J. Wolf, "Noiseless coding of correlated information sources," *IEEE Transactions on Information Theory*, vol. 19, no. 4, pp. 471–480, 1973.
- [24] A. Wyner and J. Ziv, "The rate-distortion function for source coding with side information at the decoder," *IEEE Transactions on Information Theory*, vol. 22, no. 1, pp. 1–10, 1976.
- [25] B. Girod, A.M. Aaron, S. Rane, and D. Rebollo-Monedero, "Distributed video coding," *Proceedings of the IEEE*, vol. 93, no. 1, pp. 71–83, 2005.
- [26] V.K. Goyal, A.K. Fletcher, and S. Rangan, "Compressive sampling and lossy compression: Do random measurements provide an efficient method of representing sparse signals?," *Signal Processing Magazine, IEEE*, vol. 25, no. 2, pp. 48–56, March 2008.
- [27] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [28] E. van den Berg and M. P. Friedlander, "In pursuit of a root," Tech. Rep. TR-2007-19, Department of Computer Science, University of British Columbia, June 2007, Preprint available at http://www.optimization-online.org/DB_HTML/2007/06/1708.html.
- [29] D. Varodayan, A. Aaron, and B. Girod, "Rate-adaptive codes for distributed source coding," *Signal Processing*, vol. 86, no. 11, pp. 3123–3130, 2006.
- [30] X. Artigas, J. Ascenso, M. Dalai, S. Klomp, D. Kubasov, and M. Oualet, "The DISCOVER Codec: Architecture, Techniques and Evaluation," *Picture Coding Symposium, Lisbon, Portugal, November*, vol. 6, pp. 14496–10, 2007.
- [31] T.M. Cover and J.A. Thomas, *Elements of information theory*, Wiley New York, 1991.
- [32] L. Rabiner and B.H. Juang, *Fundamentals of speech recognition*, Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1993.
- [33] E.J. Candes, M.B. Wakin, and S.P. Boyd, "Enhancing Sparsity by Reweighted ℓ_1 Minimization," *Preprint*, 2007.
- [34] R. Bernardini, M. Naccari, R. Rinaldo, M. Tagliasacchi, S. Tubaro, and P. Zontone, "Rate allocation for robust video streaming based on distributed video coding," *Signal Processing: Image Communication*, vol. 23, no. 5, pp. 391–403, 2008.
- [35] "Audacity web site," <http://audacity.sourceforge.net/>.