

Mixed Watermarking-Fingerprinting Approach for Integrity Verification of Audio Recordings

Emilia Gómez¹, Pedro Cano¹, Leandro de C. T. Gomes², Eloi Batlle¹, Madeleine Bonnet²

¹ Music Technology Group, Pompeu Fabra University, Spain

{emilia.gomez, pedro.cano, eloi.batlle}@iuia.upf.es, <http://www.iua.upf.es/mtg/>

² InfoCom-Crip5, Université René Descartes, Paris, France

{tgomes, bonnet}@math-info.univ-paris5.fr, <http://www.math-info.univ-paris5.fr/crip5/infocom/>

Abstract—We introduce a method for audio-integrity verification based on a combination of watermarking and fingerprinting. An audio fingerprint is a perceptual digest that holds content information of a recording and allows one to identify it from other recordings. Integrity verification is performed by embedding the fingerprint into the audio signal itself by means of a watermark. The original fingerprint is reconstructed from the watermark and compared with a new fingerprint extracted from the observed signal. If they are identical, the signal has not been modified; if not, the system is able to determine the approximate locations where the signal has been corrupted. The watermarked signal may undergo content preserving transformations, such as resampling or D/A and A/D conversion, without triggering the corruption alarm.

I. INTRODUCTION

IN many applications, the integrity of an audio recording must be unquestionably established before the signal can actually be used, i.e. one must be sure that the recording has not been modified without authorization.

When dealing with speech, some application contexts in which integrity must be ensured are:

- protection of previously recorded testimonies that are to be used as evidence before a court of law;
- protection of recorded interviews, which could be edited for malicious purposes.

Regarding music applications, some examples are:

- integrity verification of radio or television commercials to ensure they air as negotiated;
- integrity verification of music aired by radio stations or distributed on the Internet, in order to report unauthorized modifications.

Integrity verification systems have been proposed as an answer to this need. Two classes of methods are well suited for these applications: *watermarking*, which allows one to embed data into the signal, and *fingerprinting*, which consists in extracting a “signature” (the fingerprint) from the audio signal.

After a conceptual description of integrity verification schemes based solely on fingerprinting or watermarking, we propose a mixed approach that takes advantage of both technologies.

Music Technology Group, Institut Universitari de l'Audiovisual (IUA), Universitat Pompeu Fabra (UPF), Area Estació de França, Passeig de Circumval·lació 8, 08003 Barcelona, Spain. Phone: +34 93 542 2201, Fax: +34 93 542 2202.

InfoCom-Crip5, UFR de Mathématiques et Informatique, Université René Descartes, 45 rue des Saints-Pères, 75270 Paris cedex 06, France. Phone: +33 01 44 55 35 24, Fax: +33 01 44 55 35 35.

II. INTEGRITY VERIFICATION SYSTEMS: A CONCEPTUAL REVIEW

A. Watermarking-Based Systems

Audio watermarking (see [1] for an introductory paper on watermarking and information hiding) consists in embedding a mark (the *watermark*) into an audio signal. This mark is also an audio signal and carries data that can be retrieved from the watermarked signal. Ideally, the watermark should not introduce any perceptible degradation in the signal, which means that the original and the watermarked signals should sound exactly the same to the listener. To a certain extent, this can be achieved by using psychoacoustic models such as those found in perceptual coding [2], [3].

In watermarking-based integrity-verification systems, the integrity of a previously watermarked audio signal is determined by checking the integrity of the watermark. We define three classes of integrity-verification systems based on watermarking:

1. Methods based on fragile watermarking, which consist in embedding a fragile watermark into the audio signal (e.g. a low-power watermark). If the watermarked signal is edited, the watermark must no longer be detectable. By “edited”, we understand any modification that could corrupt the content of a recording. “Cut-and-paste” manipulations (deletion or insertion of segments of audio), for example, must render the watermark undetectable. In contrast, content-preserving manipulations (such as lossy compression with reasonable compression rates or addition of small amounts of channel noise) should not prevent watermark detection.

Extremely fragile watermarks can also be used to verify whether a signal has been manipulated in any way, even without audible distortion. For example, a recording company can watermark the content of its CDs with a very fragile watermark. If songs from this CD are bit-compressed (e.g. in MPEG format), then decompressed and recorded on a new CD, the watermark would not be detected in the new recording, even if the latter sounds exactly as the original one to the listener. A CD player can then check for the presence of this watermark; if no watermark is found, the recording has necessarily undergone illicit manipulations and the CD is refused. The main flaw in this approach is its inflexibility: as the watermark is extremely fragile, there is no margin for the rights owner to define any allowed signal manipulations (except for the exact duplication of the audio signal).

2. Methods based on semi-fragile watermarking, which are a variation of the previous class of methods. The idea consists in circumventing the excessive fragility of the watermark by increasing its power. This semi-fragile watermark is able to resist slight modifications in the audio signal but becomes undetectable when the signal is more significantly modified. The difficulty in this approach is the determination of an appropriate “robustness threshold” for each application.

3. Methods based on robust watermarking, which consist in embedding a robust watermark into the audio signal. The watermark is supposed to remain detectable in spite of any manipulations the signal may undergo. Integrity is verified by checking whether the information contained in the watermark has been corrupted or not. As the original content of the watermark must be known to the system, this content may have to be stored elsewhere.

Watermarking-based integrity-verification systems depend entirely on the reliability of the watermarking method. However, an audio signal often contains short segments that are difficult to watermark due to localized unfavorable characteristics (e.g. very low power or ill-conditioned spectral characteristics); these segments will probably lead to detection errors, particularly after lossy transformations such as resampling or MPEG compression. In integrity-verification applications, this is a serious drawback, since it may not be possible to decide reliably whether unexpected data are a consequence of intentional tampering or “normal” detection errors.

B. Fingerprinting-Based Systems

Audio fingerprinting or **content-based identification** (CBID) methods extract relevant acoustic characteristics from a piece of audio content. The result is a perceptual digest, the *fingerprint*, that acts as a kind of digital signature of the audio signal. If the fingerprints of a set of recordings are stored in a database, each of these recordings can be identified by extracting its fingerprint and searching for it in the database.

In fingerprinting-based integrity-verification systems, the integrity of an audio signal is determined by detecting changes in its fingerprint. These systems operate in three steps: (1) a fingerprint is extracted from the original audio recording, (2) this fingerprint is stored in a trustworthy database, and (3) the integrity of a recording is verified by extracting its fingerprint and comparing it with the original fingerprint stored in the database. If the transmission is digital, the fingerprint can be sent as part of a header [4].

Some fingerprinting methods evolve from traditional cryptographic hashing algorithms. An integrity-verification system can be implemented by directly applying such algorithms:

1. Methods sensitive to data modification, based on hashing methods such as MD5 [5]. This class of methods is appropriate when the audio recording is not supposed to be modified at all, since a single bit flip is sufficient for the fingerprint to change. Some robustness to slight signal modifications can be obtained by not taking into account the least-significant bits when applying the hash function.

In order to avoid sensitivity to content preserving operations, there has been an evolution towards content-based fingerprints:

2. Methods sensitive to content modification, based on fingerprinting methods that are intended to represent the content of an audio recording (such as AudioDNA [6]). This class of methods is appropriate when the integrity check is not supposed to be compromised by operations that preserve audio content (in a perceptual point of view) while modifying binary data, such as resampling, lossy compression or D/A and A/D conversion.

The main disadvantage of fingerprinting-based methods is the need of additional metadata (the original fingerprint) in the integrity-check phase. This requires access to a database or the insertion of the fingerprint in a header (not appropriate for analog audio streams) [4].

III. A COMBINED WATERMARKING-FINGERPRINTING SYSTEM

The branch of integrity-verification that combines watermarking and fingerprinting is known as *self-embedding* [4]. The idea consists in extracting the fingerprint of a signal and storing it in the signal itself through watermarking, thus avoiding the need of additional metadata during integrity check.

Some methods based on this idea have already been described in the literature, particularly for image and video [7], [8]. Shaw proposed a system [9] that embeds an encrypted hash into digital documents, also including audio. This approach inherits the limitations of hashing methods with respect to fingerprinting: hashing methods are sensitive to content preserving transformations (see section II.B).

We propose an integrity verification approach that combines a fingerprinting method representing the content of an audio recording and a robust watermarking algorithm. Fig. 1 presents a general scheme of this mixed approach.

First, the fingerprint of the original recording is extracted; this fingerprint, viewed as a sequence of bits, is then used as the information to be embedded into the signal through watermarking. As the watermark does not affect the perceptual quality of the sound, the watermarked recording should have the same fingerprint as the original recording. Thus, the integrity of this recording can be verified by extracting its fingerprint and comparing it with the original one (reconstructed from the watermark). This procedure will be detailed in the following sections.

A. Requirements

We mention below some of the requirements that are expected to be satisfied by the integrity-verification system and its components:

- the fingerprint should not be modified when transformations that preserve audio content are performed;
- the watermarking scheme must be robust to such transformations;
- the bit rate of the watermarking system must be high enough to code the fingerprint information with strong redundancy;
- the method should be suitable for use with streaming audio, as the total length of the audio file is unknown in applications such as broadcasting [10].

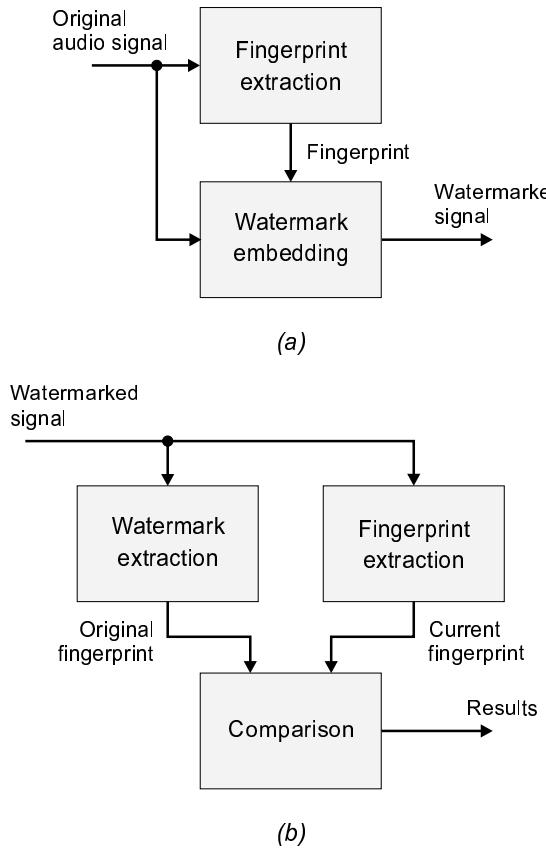


Fig. 1. Block diagram of the mixed approach for audio integrity verification: (a) embedding; (b) detection.

As shown in the following sections, the first three requirements are fulfilled by the system. The last one is also satisfied, as both the watermark and the fingerprint can be processed “on the fly”.

B. System features

The integrity-verification system can detect structural manipulations of audio signals, which correspond to the kind of tampering that must be avoided in the case of recorded testimonies or interviews. Promising results are also obtained for other distortions that perceptually affect the signal, such as:

- time stretching;
- pitch shifting;
- filtering;
- addition of noise.

The system is not only able to detect tampering, but it can also determine the approximate location where the audio signal was corrupted.

C. Implementation

C.1 Fingerprint Extraction

The employed fingerprinting scheme considers audio as a sequence of *acoustic events*. In the case of speech signals,

for example, acoustic events can be directly associated with phonemes. In music modeling, however, this association is not straightforward. The use of musical notes, for instance, would present much more complexity, as we are dealing with polyphony; furthermore, music pieces may also contain non-harmonic sounds.

The approach consists in obtaining the relevant acoustic events — called *Audio Descriptor Units* (ADU or AudioDNA) — by means of unsupervised training, i.e. without any previous knowledge of music events. The training process is performed through a modified Baum-Welch algorithm on a corpus of representative music [11].

Shortly, the system works as follows. An alphabet of representative sounds is derived from the corpus of audio signals (constructed according to the kind of signals that the system is supposed to identify). These audio units are modeled by means of Hidden Markov Models (HMM).

The audio signal is processed in a frame-by-frame analysis. A set of relevant-feature vectors is first extracted from the sound. These vectors are then normalized and sent to the decoding block, where they are submitted to statistical analysis by means of the Viterbi algorithm. The output of this chain — the fingerprint — is the most likely ADU sequence for this audio signal. This process is illustrated in Fig. 2.

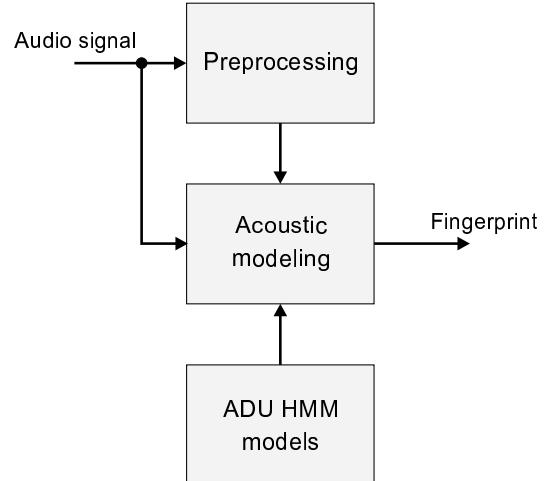


Fig. 2. Fingerprint extraction.

The resulting fingerprint is therefore a sequence of symbols (the ADUs) and time information (start time and duration). The number of different ADUs available to the system can be adjusted, as well as the output rate. The setup used in our experiments corresponds to 16 different ADUs (G_0, G_1, \dots, G_{15}) and an average output rate of 100 ADUs per minute.

C.2 Fingerprint Encoding and Watermark Embedding

Each 8-s segment of the audio signal is treated individually in order to allow for streaming-audio processing. The fingerprint is converted into a binary sequence by associating a unique four-bit pattern to each of the 16 possible ADUs; thus, the av-

verage fingerprint bit rate is approximately 7 bits/s. In our experiments, the watermark bit rate is set to 125 bits/s, allowing the fingerprint information to be coded with huge redundancy (which minimizes the probability of error during its extraction). A simple repetition code is employed, with a particular 6-bit pattern (011110) serving as a delimiter between repetitions. To avoid confusion between actual data and delimiters, every group of four or more consecutive bits “1” in the data receives an additional bit “1”, which is suppressed in the detection phase.

Fingerprint data is embedded into the audio signal by means of a watermark. The watermarking system used in our experiments is represented in Fig. 3. The analogy between watermark-

with a specific input binary pattern. The modulator output is produced by concatenating codebook vectors according to the input data sequence.

To ensure watermark inaudibility, the modulator output is spectrally shaped through filtering according to a masking threshold (obtained from a psychoacoustic model). This procedure, repeated for each window of the audio signal (≈ 10 ms), produces the watermark. The watermarked signal is obtained by adding together the original audio signal and the watermark.

As transmission and reception must be synchronized, the transmitted data sequence also carries synchronization information. This sequence is structured in such a way that detected data is syntactically correct only when the detection is properly synchronized. If synchronism is lost, it can be retrieved by systematically looking for valid data sequences. This resynchronization scheme, based on the Viterbi algorithm, is detailed in [12] and [13].

C.3 Watermark Detection and Fingerprint Decoding

For each window of the received signal, the watermark signal is strengthened through Wiener-filtering and correlation measures with each codebook entry are calculated. The binary pattern associated with the codebook entry that maximizes the correlation measure is selected as the received data. The syntactic consistency of the data is constantly analyzed to ensure synchronization, as described in the previous section.

The output binary sequence is then converted back into ADUs. For each 8-s audio segment, the corresponding fingerprint data is repeated several times in the watermark (16 times in average). Possible detection errors (including most errors caused by malicious attacks) can then be corrected by a simple majority rule, providing a replica of the original fingerprint of the signal.

C.4 Matching and Report

Finally, the fingerprint of the watermarked signal is extracted (through the same procedure presented in section III-C.1) and compared with the original fingerprint obtained from the watermark. If the two sequences of ADUs match perfectly, the system concludes that the signal has not been modified after watermarking; otherwise, the system determines the instants associated to the non-matching ADUs, which correspond the approximate locations where the signal has been corrupted. Identical ADUs slightly shifted in time are considered to match, since such shifts may occur when the signal is submitted to content-preserving transformations.

IV. SIMULATIONS

A. Experimental conditions

Results of cut-and-paste tests are presented for four 8-s test signals: two songs with voice and instruments (signal “cher”, from Cher’s “Believe”, and signal “estrella_morente”, a piece of flamenco music), one song with voice only (signal

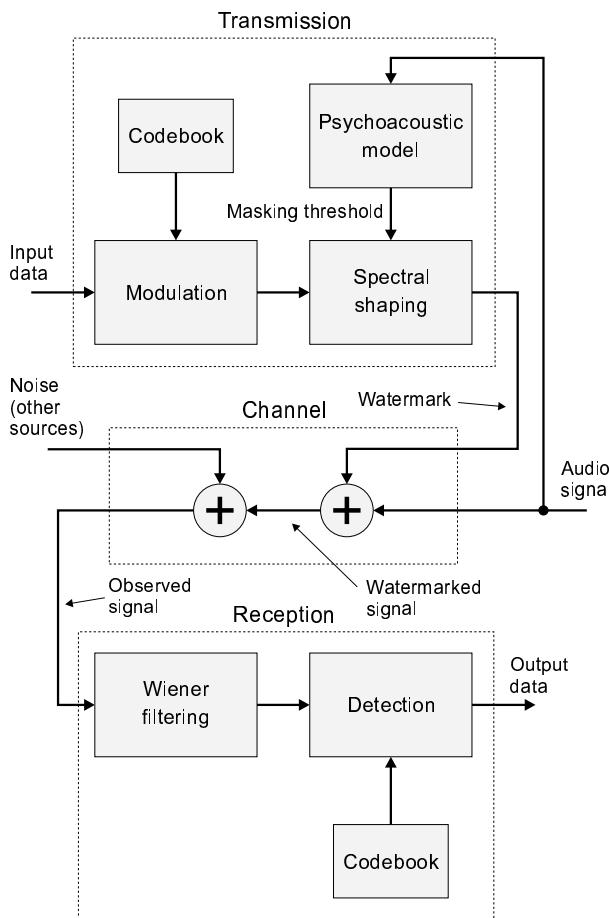


Fig. 3. Watermarking system.

ing and digital communications is emphasized in the figure: watermark synthesis corresponds to transmission (with the watermark as the information-bearing signal), watermark embedding corresponds to channel propagation (with the audio signal as channel noise), and watermark detection corresponds to reception.

The watermark signal is synthesized from the input data by a modulator. In order to obtain a watermark that is spread in frequency (so as to maximize its power and increase its robustness), a codebook containing white, orthogonal Gaussian vectors is used in the modulator. The number of vectors is a function of the desired bit rate. Each codebook entry is associated

“svega”, Suzanne Vega’s “Tom’s diner”, a cappella version), and one speech signal (signal “the_breakup”, Art Garfunkel’s “The breakup”). The signals were sampled at 32 kHz and were inaudibly watermarked with a signal to watermark power ratio of 23 dB in average.

B. Results

Fig. 4 shows the simulation results for all test signals. For each signal, the two horizontal bars represent the original signal (upper bar) and the watermarked and attacked signal (lower bar). Time is indicated in seconds on top of the graph. The dark-gray zones correspond to attacks: in the upper bar, they represent segments that have been *inserted* into the audio signal, whereas in the lower bar they represent segments that have been *deleted* from the audio signal. Fingerprint information (i.e. the ADUs) is marked over each bar.

For all signals, the original fingerprint was successfully reconstructed from the watermark. Detection errors introduced by the cut-and-paste attacks were eliminated by exploiting the redundancy of the information stored in the watermark.

A visual inspection of the graphs in Fig. 4 shows that the ADUs in the vicinities of the attacked portions of the signal were always modified. These corrupted ADUs allow the system to determine the instant of each attack within a margin of approximately ± 1 second.

For the last signal (“the_breakup”), we also observe that the attacks induced two changes in relatively distant ADUs (approximately 2 s after the first attack and 2 s before the second one). This can be considered a false alarm, since the signal was not modified in that zone.

V. ADVANTAGES OF THE MIXED APPROACH

In this section, we summarize the main advantages of the mixed approach in comparison with other integrity-verification methods:

- No side information is required for the integrity test; all the information needed is contained in the watermark or obtained from the audio signal itself. This is not the case for systems based solely on fingerprinting, since the original fingerprint is necessary during the integrity test. Systems based solely on watermarking may also require side information, as the data embedded into the signal cannot be deduced from the signal itself and must be stored elsewhere;
- Slight content-preserving distortions do not lead the system to “false alarms”, since the fingerprint and the watermark are not affected by these transformations. Hashing methods (such as MD5) and fragile watermarks generally do not resist such transformations;
- In general, localized modifications in the audio signal also have a localized effect on the fingerprint, which enables the system to determine the approximate locations where the signal has been corrupted. This is not the case for simple hashing methods, since the effect of localized modifications may propagate to the entire signal;
- Global signal modifications can also be detected by the system; in this case, the entire fingerprint will be modified and/or the watermark will not be successfully detected;

- This method is well suited for streaming audio, since all the processing can be done in real time.

VI. CONCLUSIONS

In this paper, we have presented a system for integrity verification of audio recordings based on a combination of watermarking and fingerprinting. By exploiting both techniques, our system avoids most drawbacks of traditional integrity-verification systems based solely on fingerprinting or watermarking. Unlike most traditional approaches, no side information is required for integrity verification. Additionally, the effect of localized modifications generally do not spread to the rest of the signal, enabling the system to determine the approximate location of such modifications. Experimental results confirm the effectiveness of the system.

As next steps in this research, we will consider possible developments in order to further increase overall system reliability, particularly in what concerns false alarms (i.e. signal modifications detected after content-preserving transformations or in zones where the signal was not modified). More efficient coding schemes will also be considered for fingerprint encoding prior to embedding. A more exhaustive evaluation of the system to examine the performance with respect to a taxonomy of different modifications is needed.

ACKNOWLEDGMENTS

Leandro de C. T. Gomes would like to thank CNPq for financial support.

REFERENCES

- [1] F. A. P. Petitcolas, R. J. Anderson, and M. G. Kuhn, *Information hiding — a survey*, Proceedings of the IEEE, special issue on protection of multimedia content, 87(7):1062-1078, July 1999.
- [2] L. Boney, A. Tewfik, and K. Hamdy, *Digital watermarks for audio signals*, International Conference on Multimedia Computing and Systems, Hiroshima, June 1996.
- [3] M. Perreau Guimarães, *Optimisation de l’allocation des ressources binaires et modélisation psychoacoustique pour le codage audio*, PhD Thesis, Université Paris V, Paris, 1998.
- [4] C.-P. Wu and C.-C. Jay Kuo, *Speech content integrity verification integrated with ITU G.723.1 speech coding*, IEEE International Conference on Information Technology: Coding and Computing, pp. 680-684, Las Vegas, April 2001.
- [5] R. Rivest, *The MD5 Message-Digest Algorithm*, <<http://theory.lcs.mit.edu/rivest/Rivest-MD5.txt>>, April 1992.
- [6] Recognition and Analysis of Audio, <<http://raa.joanneum.at/>>.
- [7] J. Dittmann, A. Steinmetz, and R. Steinmetz, *Content-based digital signature for motion pictures authentication and content-fragile watermarking*, International Conference on Multimedia Computing and Systems, Florence, June 1999.
- [8] J. Dittmann, *Content-fragile watermarking for image authentication*, Proceedings of SPIE, vol. 4314, Bellingham, 2001.
- [9] G. Shaw, *Digital document integrity*, 8th ACM Multimedia Conference, Los Angeles, November 2000.
- [10] R. Gennaro and P. Rohatgi, *How to sign digital streams*, Advances in Cryptology, CRYPTO’97, 1997.
- [11] E. Battle and P. Cano, *Automatic segmentation for music classification using competitive hidden markov models*, International Symposium on Music Information Retrieval, 2000.
- [12] E. Gómez, *Tatouage de signaux de musique (méthodes de synchronisation)*, DEA ATIAM thesis, ENST (Paris Télécom)-IRCAM (Centre Georges Pompidou), 2000, <<http://www.iua.upf.es/mtg/>>.
- [13] L. de C. T. Gomes, E. Gómez, and N. Moreau, *Resynchronization methods for audio watermarking*, 111th AES Convention, New York, November 2001, <<http://www.iua.upf.es/mtg/>>.

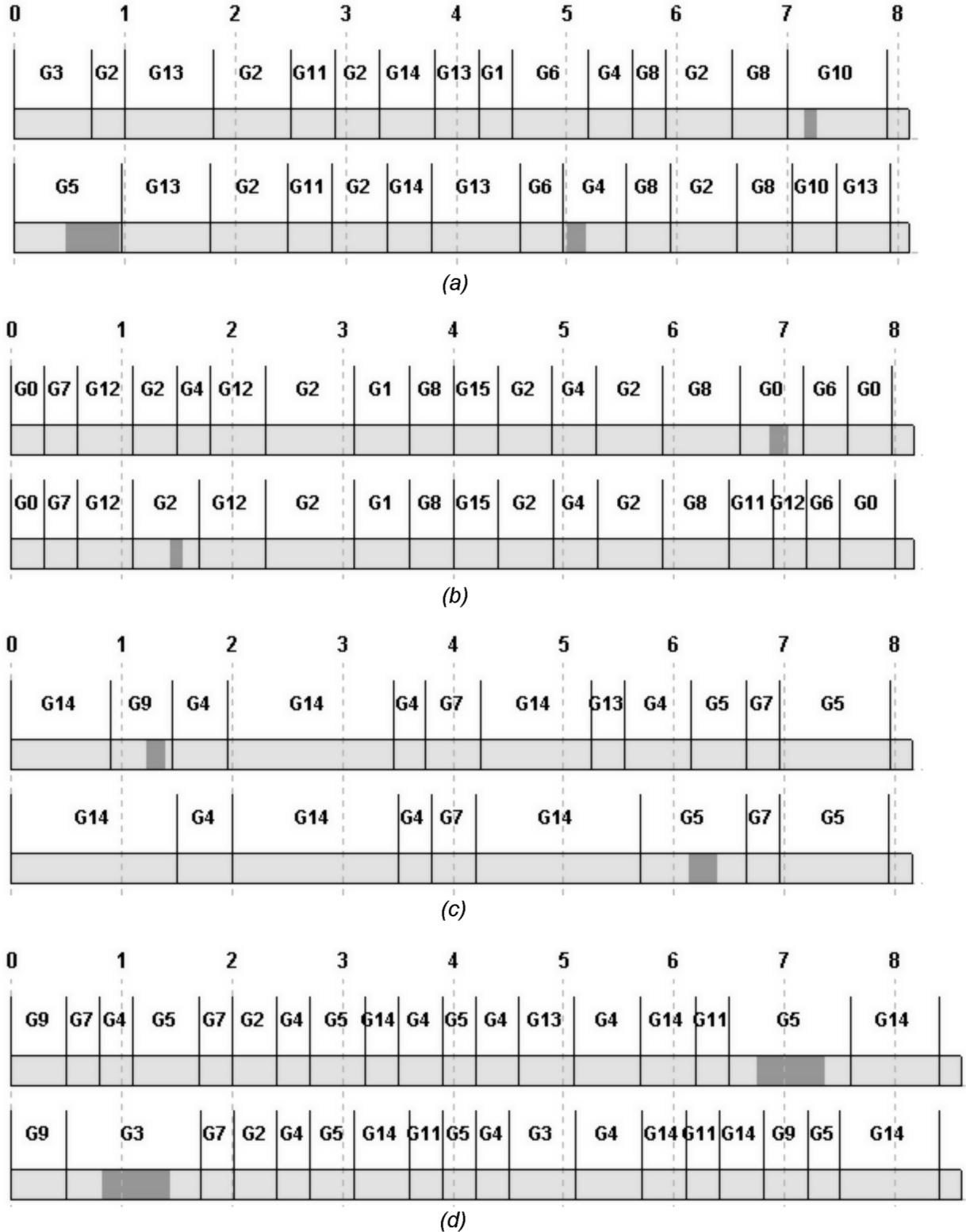


Fig. 4. Simulation results: (a) signal “cher”; (b) signal “estrella_morente”; (c) signal “svega”; (d) signal “the_breakup”.