

Research on Security Transmission of Perceptual Hash Values Based on ECC and Digital Watermarking

Zhang Qiuyu¹, Xing Pengfei¹, Liu Yangwei¹, Zhang Qianyun² and Huang Yibo¹

¹ School of Computer and Communication, Lanzhou University of Technology,
Lanzhou, 730050, China

² School of Communication & Information Engineering, Shanghai University,
Shanghai, 200444, China

zhangqylz@163.com, xingpengfei0202@126.com, cc1000cc@126.com

Abstract

In this paper, a secure method of transmitting perceptual hash is proposed based on error correcting codes (ECC) and digital watermarking, aiming at the fact that perceptual hash string used in audio authentication is easy to alter after been attacked in transmission and the need of extra channel. In this method we embed the binary perceptual hash values into the compressed audio as a digital watermark so that authentication data is dispersed. An ECC is used to pre-process perceptual hash sequence to ensure that speech signal with watermarking can be extracted watermarking information (the perceptual hash values) after been attacked. The experimental result illustrate that this method embeds information without influencing on the content. It is robust to noise and can prevent from common attack. At the same time, this method ensures the hearing transparency of audio perceptual content, the high efficiency of watermark extraction and the efficiency and security of algorithm.

Keywords: Audio perceptual hash, Digital watermarking, Error correcting code (ECC), Perceptual hash values, Security analysis

1. Introduction

At present, the research of perceptual hashing [1] mainly consists of two aspects, identification and authentication. With the appearance of many perceptual hash algorithms, researchers started to discuss problems in performance evaluation of perceptual hash, which are mainly about basic properties of perceptual hash, such as robustness, discrimination, compactness and key dependence [2, 3] rather than security of which research is in the initial stage [4].

The security of authentication system consists of two aspects, the security of authentication algorithm and the security of authentication information. Currently the algorithm that most of the literatures use to evaluate security is: whether the perceptual hash value changes randomly when few numbers of key changes. The perceptual hash of image or audio ought to be difficult to be forged and estimated when the key is unknown. Hence, what many algorithms focus on is how to embed efficient key during perceptual features extracting or hash modeling rather than the security of authentication information, *i.e.*, the security of perceptual hash values.

There is a secure channel transferring authentication information exists in perceptual hash matching framework which increases transmission cost. The usual handling that adds the perceptual hash to the end of digital representation of authentication object locates the perceptual hash values in an area that easy to be attacked. Considering that perceptual hash

just reflects perceptual content of multimedia object, and does not be influenced by the change of digital representation of multimedia object, hence, we do not need to consider the robustness and security problems resulted from the fact that hash bit sequence is hidden in the carrier. An algorithm of compressed domain watermarking is proposed in [5]. The watermark consists of synchronous codes which are revising from MDCT spectrum entropy. This method is robust to usual attack and can correctly extract the watermark in the receiver. Hu, *et al.*, [6] proposed a perceptual hash algorithm based on Encrypted domain to solve security issues in the practical application environment, and his research is based on key also. Hadmi, *et al.*, [7] proposed a safety perceptual hash system based on analysis of quantitative step which added an encryption compression module in quantization process and analyzed the security of system through Gaussian noise, JPEG compression and low-pass filtering.

To sum up, this paper proposes a secure transport method of perceptual hash values based on digital watermarking and error correcting codes (ECC), aiming at the problems in security analysis and research ideas. The proposed method can ensure the security of authentication information over an unsafe channel. The algorithm that based on MDCT spectrum entropy is firstly tested in experiments. The results are used to analyze the relationship between algorithm security and indicators. With the experimental results we come to the conclusion that the current security of audio perceptual hash is based on the secret key and special optimization for partial content tampering when constructing a perceptual hash algorithm.

2. Related Works

2.1. Digital Watermarking

As a main method of multimedia authentication, digital watermarking technology is a hot research topic in recent years. Digital watermarking is mainly divided into the robust watermark [8], semi fragile watermarking [9] and complete fragile watermarking [10]. The research starts from robustness, perceptual and embedding capacity. Watermark is used to identification, tracking, ensure the legitimate ownership using authorized access, prevent illegal copying and solves the problems in copyright and security [11].

Application of digital watermarking technology in the field of information security promotes the development of digital audio watermarking. Abundant research achievements in digital audio watermarking algorithm mainly include the Least Significant Bit method, echo hiding method and transform domain watermarking algorithms such as DFT, DCT and DWT. In [12] firstly proposed a method that embedded the watermark directly into the MPEG audio stream, respectively selected scaling factor factors and sample data of MPEG code as the embedding position. It is worthwhile to note that the direct embedding in the scale factor does not need decoding process, but will cause sound distortion in pure speech audio without background and influence perceptual content.

2.2. Error Correcting Code (ECC)

The carrier with watermark may have error code that caused by filtering, compression/decompression, cutting and other operations during transmission. In order to find transmission error and correct it, the necessary judgment data are added before transmission and the reliability of watermark carrier is improved.

Error correcting codes is a very mature technology in communication system and widely applied in practice. Because of the similarities between watermark channel and communication channel, many researchers study for the purpose of improvement in imperceptibility and robustness. The research works of both theory and experiment are about

error correcting codes function in the digital watermark system [13], including repeated code, Hamming code, BCH codes, convolution codes, Turbo codes etc.

Research of error correcting code in digital watermarking are mainly aimed at the image carriers rather than audio carriers. This paper introduces a method of convolution code to pre-process extracted perceptual hash values.

Parameter of (n, k, m) is used to describe the convolution code, where k is the number of input bits to the convolution encoder each time and n is number of correspond output bits; m represents the storage capacity and series of k -tuple; $k = m + 1$ is called constraint length. The generated n -tuple convolution codes is not only related to the current input k -tuple but also associated with previous $m - 1$ k -tuples. The number of correlative codes during coding process is $n \times m$. The error correcting performance of convolution codes increases with m and the error rate will exponentially decreases with the increase of n .

Principle of encoding and decoding of convolution codes is complex, but the process can be simplified in the MATLAB platform and the efficiency of the algorithm is ensured.

2.3. Logistic Mapping

One-dimensional Logistic map is widely used in the secure communications field. Logistic sequence has the following characteristics:

- 1) Simple form, which is easy to generate and replicate in the case that parameters and initial conditions have been determined;
- 2) Sensitive to the initial conditions. Power system's trajectories are completely different while slight differences;
- 3) Excellent statistical properties of white noise;
- 4) Ergodicity. Value of the power system, which may occur anywhere, is difficult to be predicted;
- 5) Irreversibility. We can't speculate system's state from part of value of power system, so that we can keep its security. Based Logistic map has excellent characteristics, which is in line with key requirements in the structure process of perceptual hash. This paper introduces it to ensure security of the algorithm.

The mathematical expressions of Logistic map are described as follows:

$$x_{n+1} = \mu x_n (1 - x_n) \quad (1)$$

where $\mu \in [0, 4]$ in Eq. (1), which is called bifurcation parameter, $x \in [0, 1]$. When μ values in $(3.5699456, 4]$ is especially closer to 4, the sequence iteratively generated is in pseudo-random distribution state.

3. The Proposed Scheme

3.1. Algorithmic Process

The proposed algorithm is mainly divided into three parts:

- 1) extracting hash values;
- 2) embedding hash values into the compressed speech as a watermark;
- 3) extracting the watermark and match it with the template.

The process of algorithm is shown in Figure 1.

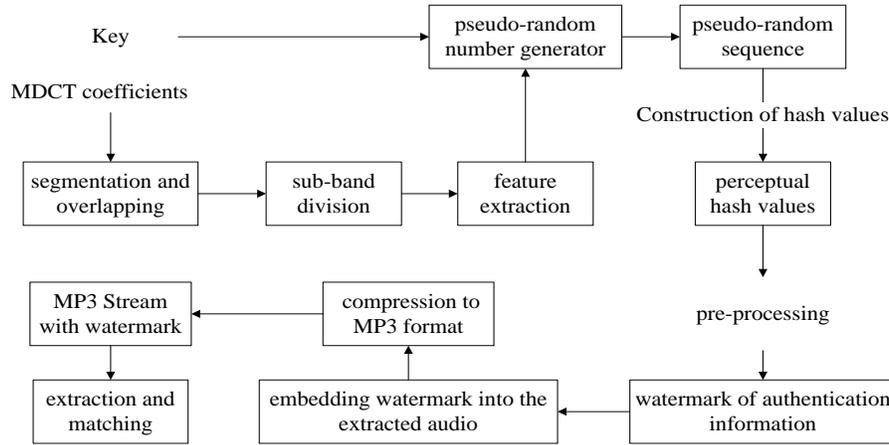


Figure 1. The Block Diagram of the Proposed Algorithm

3.2. Perceptual Hash Values Extraction and Secret Key

Strictly speaking, the speech signal is non-stationary random process. After some preprocessing techniques it can be approximated as a stationary random process, which is currently the general method of audio signal processing. Such as the original domain speech perceptual hash use multi-windowed, sub-frame and aliasing methods to make the speech signal believed to be stable signal in each 10-30ms fragment. Because MP3 files having a fixed frame structure, the section is divided into N sub-bands, which is 50% overlap between adjacent sub-bands, in order to achieve a similar preprocessing effect of sub-frame aliasing, shown in Figure 2.

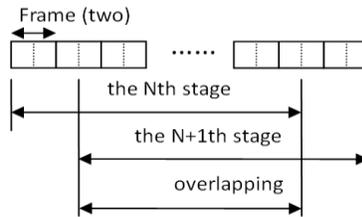


Figure 2. Sub-band Division Schematic Diagram

As shown in Figure 2, the longer the length of two segments, the higher the robustness of the algorithm, but at the same time, the lower the accuracy of the algorithm. According to the Kalker theory, the 3s audio section can fully be characterized with the binary sensing hash value, which is 256 bit. Therefore, according to the $N = \text{int}(L/123.5)$, we can calculate the number of segments contained in each section, where the L is the total number of MDCT coefficients extracted from speech clips section. For calculating easy, N takes even.

After preprocessing in above operation, calculation MDCT frequency domain, and compute the i -th sub-band energy SBE_{ij} of the j -th section. In this case, the role of the MDCT coefficient in each sub-band is similar to Fourier coefficients of the original signal power spectrum, thus SBE_{ij} can be calculated accord to Eq. (2):

$$SBE_{ij} = \sum_{m=\frac{N(i-1)}{2}+1}^{\frac{N(i+1)}{2}} \sum_{n=1}^{32} |G(m,n)|^2, \quad i \in [1,256], m, n \in Z \quad (2)$$

where N represents divided sub-band Number of fragments which has not been identified, $G(i, n)$ represents the n -th MDCT coefficients of the i -th section.

In order to obtain class probability mass function and meet their basic characteristics of statistical probability, that is to say the sum of all the elements is 1. According to Eq. (3), the MDCT spectral energy of each section divided by the total energy of sub-band, using p_{ij} said, which is defined as follows:

$$p_{ij} = \frac{SBE_{ij}}{\sum_{j=1}^N SBE_{ij}} \quad i = 1,2,3 \dots,256 \quad (3)$$

According to Eq. (4) to calculate MDCT spectral entropy of the i -th sub-band:

$$H(i) = - \sum_{j=1}^N p_{ij} \log_2 p_{ij} \quad i = 1,2,3, \dots,256 \quad (4)$$

As Eq. (5), use Logistic mapping rules to scramble, where μ is the key, and is taken close to a value of 4.

$$H_{i+1} = \mu H_i (1 - H_i) \quad (5)$$

To make the match calculated of hash value simple, at the same time, to further increase its robustness, in this paper, we utilize size comparison method of the adjacent sub-band energy spectrum entropy to form the final relatively binary perceptual hash value sequence, as follows:

$$ph(i) = \begin{cases} 0, & H(i) < H(i+1) \\ 1, & H(i) \geq H(i+1) \end{cases} \quad (6)$$

3.3. The Generation and the Embedding Process of Watermark

To encode the (2,1,2) convolutional code of the perceptual hash values generated in section 3.2. Using the grid chart can be more intuitive in showing the encoding process, which is shown in Figure 3.

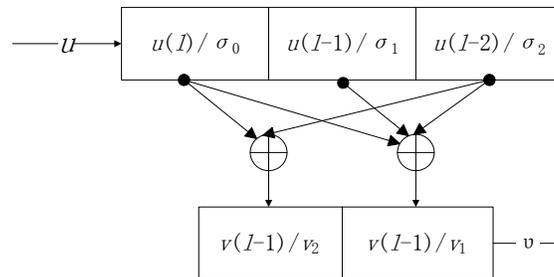


Figure 3. Convolutional Code's Encoding Process

Treat the encoded hash values as authentication information of watermark, the embedded concrete steps are as follows:

Step1: Selecting the appropriate wavelet base for the three-level wavelet decomposition of the audio signal X , Get the details of the component under different resolution levels (high frequency) d_j ($j = 1, 2, 3$) and approximation component c_3 (low frequency component).

Step2: Selecting three scale low frequency coefficient c_3 and high frequency coefficient d_3 , then framing the select the wavelet coefficients.

Step3: Seeking the maximum absolute value of wavelet coefficient per frame, as follows:

$$m_i = \max_{j=i}^N |X_{ij}| \quad (7)$$

Step4: Embedding according to Eq. (8).

$$X_i = X_i + 0.02 m_i W_i \quad (8)$$

Step5: Doing the wavelet inverse transformation, we get audio signal containing the watermark.

Step6: Recompression the PCM data for MP3 files containing the watermark information.

3.4. The Extraction and Matching of Watermark

In contrast to the process of embedded watermark, extraction operation is conducting in the process of decoding the MP3 audio file decoding into PCM code stream. This can achieve playing MP3 files and extracting the watermark at the same time. The watermark extraction algorithm time is shorter, so it will not affect the play of MP3 files.

First, the MP3 audio carriers that contain watermark are operated with attacks such as sampling, low-pass filtering, noise addition, signal enhancement, random cutting and replacement; then watermark from the signal carriers are extracted and corrected according to Viterbi decision method. The output perceptual hash values are achieved and compared with templates from database; afterwards the error rate between them is calculated.

4. Experimental Results and Analysis

4.1. Experimental Environment

In this paper, we present a full procedure of performance tests and their results. The database of 400 speech clips in our experiment, including clips with different content of Chinese and English and same content read by different people, is shown in Table 1. Each clip is compressed into MP3 format and lasts 4s. Experiment process works on the platform of MATLAB 2010b and Libmad encoder is used in re-process.

Table 1. Speech Clips

Sampling Rate	Bit Depth	Channel	Bit Rate
44100Hz	16 bits	mono	64kbps

4.2. Algorithmic Performance Analysis

Hash values are extracted from 400 speech clips within the database and then compared in pairs. The 79,800 BER results are shown in Figure 4, where the comparison of the distribution of BERs and the normal distribution is illustrated.

The Figure 4 shows that BER between speeches and with different content also has a generally normal distribution with a key. The probability distribution parameters are a mean value of $\mu=0.4869$ and standard deviation of $\sigma=0.0351$, both calculated upon the MATLAB platform. Through calculation we arrive at a low $FAR=2.3833e-18$ at a threshold $\tau=0.18$, revealing that more than 2 clips were falsely claimed as similar to entire 1018 clips in Eq. (9). Experimental results prove that the proposed algorithm meets the accuracy demands of speech identification in practical applications.

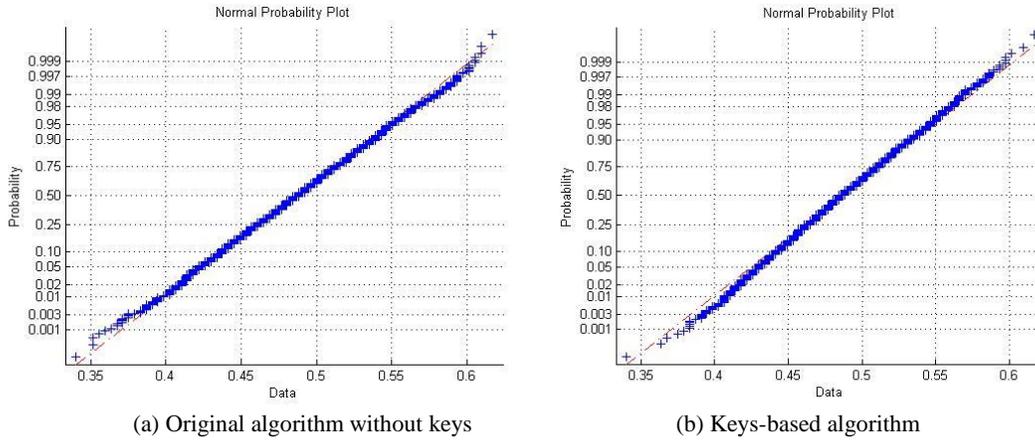


Figure 4. Comparison of Algorithm Discrimination Abilities

$$FAR(\tau) = \int_{-\infty}^{\tau} f(\alpha | \mu, \sigma) d\alpha = \int_{-\infty}^{\tau} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(\alpha - \mu)^2}{2\sigma^2}} d\alpha \quad (9)$$

All of the 100 MP3 speech clips are subjected to the following procedures:

- Increase the volume by 50%.
- Reduce the volume by 50%.
- Resample consisting of subsequent down and up sampling to 22.05 kHz and 44.10 kHz.
- Echo addition with attenuation of 60%, time delay of 300ms and initial strength of 20% and 10%
- Noise addition with center frequency of 0~4 kHz
- Low-pass filtering, using a fifth order Butterworth filter with cut-off frequency of 2 kHz.

Each of the operations can preserve the perceptual content of speech signals except for the last signal. Hash values are extracted from speech clips processed with the first five content-preserving operations and the BER between the hash values are determined. The values are extracted from clips with the same perceptual content. The resulting bit error rates are shown in Figure 5 (with same perceptual content) and Table 2.

Table 2. Average BER

Operations	Without Key	Key-based
Volume down	0.0096	0.1002
Volume up	0.0179	0.0174
Echo addition	0.1872	0.1881
Resample	0.0068	0.0074
Noise addition	0.0415	0.0412
Low-pass filtering	0.2746	0.2752

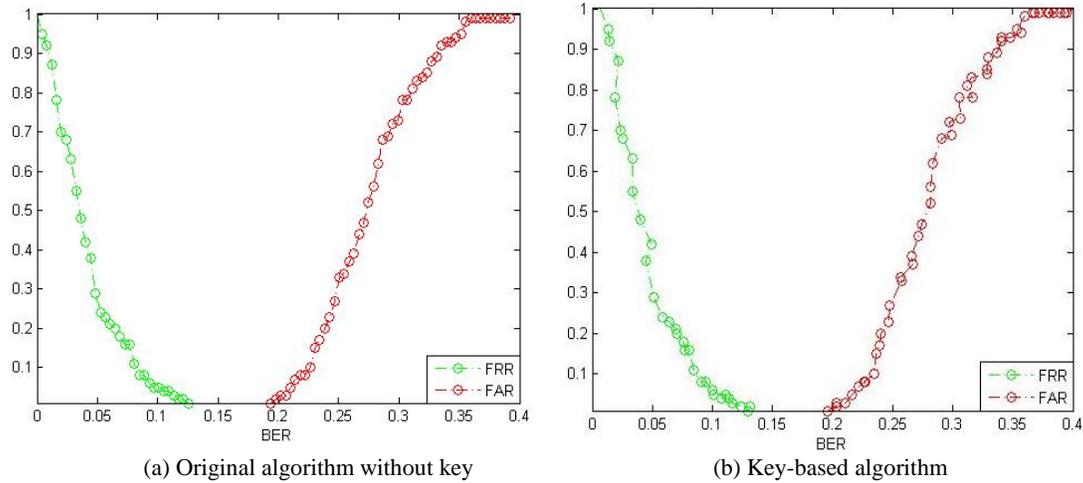


Figure 5. Comparison of Algorithm Robustness

Afterwards, FAR curve of speech clips are subjected to low-pass filtering that is drawn from within the same FRR curve, coordinate system. The interval of discrimination between 0.14 and 0.19 makes it possible for the proposed algorithm to certify clips performed by content-preserving operations.

It can be seen from this pair of figures that there are hardly any differences in the interval of discrimination and the performance of algorithm is maintained.

The time consumption is bound to increase due to the randomization mechanism of Logistic mapping. Experiments are repeated 20 times using 100 groups of extracted MDCT coefficients. The extraction time of hash values is increased by 0.03s, which means the efficiency of algorithm is not influenced randomization mechanism.

4.3. Security Analysis of Algorithm

As described in [11], we can know that security of the perceptual hash can be reflected through scrambling and diffusion ability. In the security test, this paper firstly tests the key's ability of scrambling. We can see from Figure 6 that when the key takes different values, all the bit error rate between $ph=PH(I, k)$ and $ph'=PH(I, k')$ are centrally distributed in about 0.5.

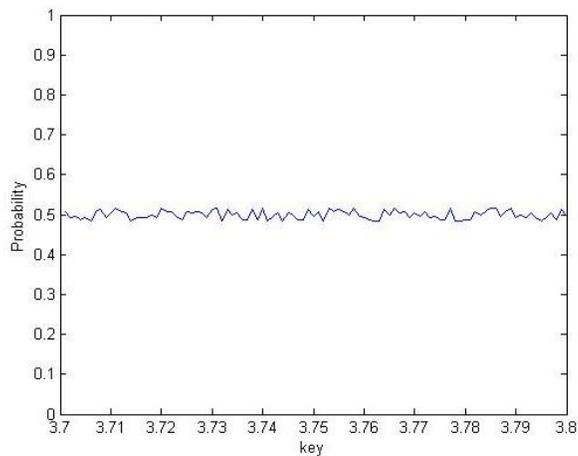


Figure 6. Keys Scrambling Capability

Diffusion ability test are mainly conducted through attacking using partial replacement. With the rise of substitution's degree, the bit error rate between $ph=PH(I, k)$ and $ph'=PH(I, k')$ is also rising but the rising rate is relatively slow, only when the degree of substitution is close to 100%, the matching results of BER is close to 0.5. In [14] also puts out that in diffusion performance test, all the measuring method used in the experiment showed the very bad diffusion ability, shows the same robustness of algorithm for other attack types. Through the test of a few more classical algorithm, in [10] got the similar conclusions, which is consistent with the experimental results in this paper.

Combining the proposed algorithm with the experimental results in [14], we can draw to the following conclusions for the algorithm's security of the Speech perceptual hashing:

- 1) Good scrambling performance of existing algorithms is mainly provided by the pseudo-random sequence generator;
- 2) For the security applications environment of authenticity of the speech perceptual hash's part content (*e.g.*, tamper detection, *etc.*,) the current perceptual of hash algorithm does not provide enough sensitive diffusion ability for part changes in speech.
- 3) Sensitive diffusion ability for part changes in speech need algorithm to be targeted optimized for one part content attack. If an algorithm can meet the requirements of scrambling and diffusion properties at the same time, then the algorithm is safe in theory.

4.4. Security of Hash Transmission

First, the MP3 audio carriers that contain watermark are operated with attacks such as sampling, low-pass filtering, noise addition, signal enhancement, random cutting and replacement; then watermark from the signal carriers are extracted and corrected according to Viterbi decision method. The output perceptual hash values are achieved and compared with templates from database; afterwards the error rate between them is calculated. The experimental results are shown in Table 3.

Table 3. Average BER

Operation	BER
Echo Addition	0.0327
Resample	0.0052
Gaussian white noise	0.0108
Low-pass filtering	0.1325

It can be seen from the data in Table 3 that hash values extracted from the operated audio and the original audio exist subtle error rate. When a grayscale image is used as authentication information, this subtle error rate will not affect perceptual content of the authentication information. But perceptual hash values have no subjective perceptibility when been used as authentication information; therefore the BER may affect the final results in case of small matching threshold. Thus, the proposed algorithm can provide a relatively safe transmission scheme, but the details still need to be further optimization in-depth.

4.5. Analysis of Auditory Perceptual

One of the watermarking system requirements is the imperceptibility, which is also consistent with the requirements of auditory perceptual. The evaluation process of the watermarked speech signal quality is similar to evaluation of the quality audio signal processed by the audio encoding decoder, including the evaluation of sound quality and distortion. Objective and subjective evaluation are common evaluation methods.

Signal noise ratio (*SNR*) is one of the most commonly used distortion measurement in time domain that can be used as an objective evaluation method to determine the extent of the carrier digital change. The calculation of *SNR* is shown in Eq. (10):

$$SNR = 10 \lg \frac{\sum_{i=0}^{N-1} x(i)^2}{\sum_{i=0}^{N-1} [x(i) - x'(i)]^2} \quad (10)$$

where $x'(i)$ and $x(i)$ respectively the original speech signal and the watermarked speech signal.

The embedded watermark changes the digital representation of the carrier and the difference of carriers after embedded process shown in the waveform of Figure 7.

As can be seen in Figure 7, the difference between waveforms is not obvious after been embedded. The *SNR* value is 37.69 dB between the original speech and the watermarked signal, which indicates that the difference is barely felt.

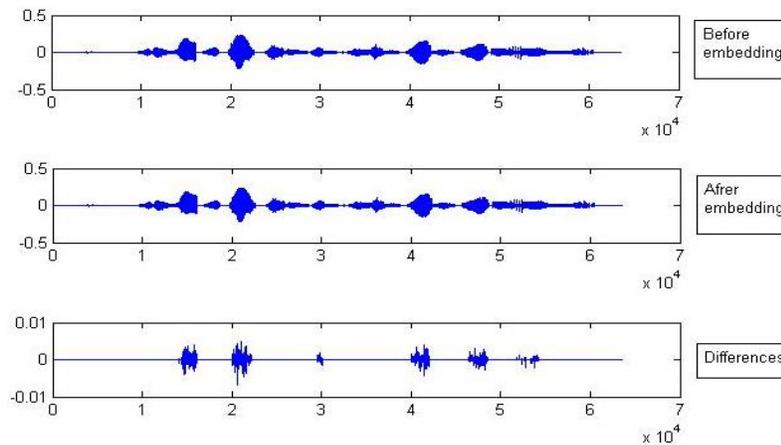


Figure 7. Waveform Graph

Objective Diff-grade (SDG) is used to in the subjective listening test. The original and the embed speeches are listened by 10 people, the resulting mean value 0.1 are classified and calculated according to SDG standard, which is approximately to zero, The result means that the watermarked and the original speech are difficult to distinguish and transparency of the watermark is maintained in the proposed algorithm.

5. Conclusions

Aiming at the less of method of security analysis and evaluation criterion in research of audio perceptual hash, a secure method transmitting perceptual hash based on error correcting code and digital watermarking is proposed in this paper. This method uses the generated hash values (authentication information) to conduct (2,1,2) convolutional encoding to embed watermark though Logistic mapping. The results show that there is no loss in performance for algorithm controlled by key, carrier signals can extract authentication information in channel after been attacked usually, so proving that the algorithm proposed in this paper can ensure the security of perceptual hash transmission even if there is no security of channel.

Nowadays, the security of audio perceptual hash algorithm is lack of scientific and specific definition, also the efficient and widely accepted evaluation indices. If extend the security of algorithm as scrambling and diffusion, the test of diffusion will mainly depend on an

optimizing test for a certain local attack without universality. So, ensuring security needs of authentication application, building evaluation indices and estimate model for security, are problems need to be solved in follow-up research.

Acknowledgments

This work is supported by the National Natural Science Foundation of China (No. 61363078), the Natural Science Foundation of Gansu Province of China (No. 1212RJZA006, No. 1310RJYA004). The authors would like to thank the anonymous reviewers for their helpful comments and suggestions.

References

- [1] Y. H. Jiao, "Research on Perceptual Audio Hashing", Harbin Institute of Technology, Heilongjiang, China, (2010).
- [2] N. Wenyin, J. X. Wang and S. P. Chen, "The Analysis of Key Technology of the Multimedia Data Content Detection System Based on Perceptual Hash", *Sensor Letters*, vol. 11, no. 4, (2013), pp. 715-718.
- [3] N. Chen and H. D. Xiao, "Perceptual audio hashing algorithm based on Zernike moment and maximum-likelihood watermark detection", *Digital Signal Processing*, vol. 23, no. 4, (2013), pp. 1216-1227.
- [4] O. Koval, S. Voloshynovskiy and D. Beekhof, "Security analysis of robust perceptual hashing", *Proc. SPIE 6819, Security, Forensics, Steganography, and Watermarking of Multimedia Contents X*, 681906, (2008), pp. 1-10.
- [5] T. Liu, Z. F. Ma and M. Jiang, "Robust MP3 Audio Watermarking Algorithm Based on MDCT Spectral Entropy Recognition", *Computer Science*, vol. 38, no. 12, (2012), pp. 113-117.
- [6] D. H. Hu, B. Su and S. L. Zheng, "Secure Architecture and Protocols for Robust Perceptual Hashing", *Proc. of the 9th IEEE International Conference on Computational Intelligence and Security (CIS)*, Leshan, Sichuan, China, (2013), pp. 550-554.
- [7] A. Hadmi, W. Puech and B. Ait Es Said, "A robust and secure perceptual hashing system based on a quantization step analysis", *Signal Processing: Image Communication*, vol. 28, no. 8, (2013), pp. 929-948.
- [8] G. Zhu and J. L. Zhang, "Adaptive robust watermarking algorithm based on SVD and wavelet packet transform", *Application Research of Computers*, vol. 30, no. 4, (2013), pp. 1230-1233.
- [9] H. X. Wang and M. Q. Fan, "Centroid-based Semi-fragile Audio Watermarking in Hybrid Domain", *Science in China Series F-Information Sciences*, vol. 53, no. 3, (2010), pp. 619-633.
- [10] Y. R. Huo, H. J. He and F. Chen, "Fragile Watermarking Algorithm for JPEG Images Based on Neighborhood Comparison and its Performance Analysis", *Journal of Software*, vol. 23, no. 9, (2012), pp. 2510-2521.
- [11] J. Panda and M. Kumar, "Application of energy efficient watermark on audio signal for authentication", *Proc. of the IEEE International Conference on Computational Intelligence and Communication Networks (CICN)*, Gwalior, India, (2011), pp. 202-206.
- [12] S. Jiang, Z. Lu and H. Wu, "DCT based multipurpose watermarking technique for image copyright notification and protection", *Journal of Harbin Institute of Technology*, vol. 11, no. 3, (2004), pp. 237-239.
- [13] Y. H. Ma and C. Y. Xu, "Survey and realization on information hiding based on error-correcting code", *Application Research of Computers*, vol. 29, no. 7, (2012), pp. 2686-2689.
- [14] H. Zhang, "Research on Benchmark and Algorithm of Image Perceptual Hashing", Harbin Institute of Technology, Heilongjiang, China, (2009).

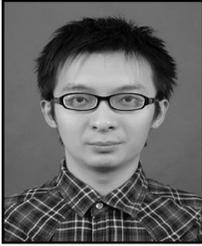
Authors



Zhang Qiuyu, Researcher/PhD supervisor, graduated from Gansu university of technology in 1986, and then worked at school of computer and communication in Lanzhou university of technology. He is vice dean of Gansu manufacturing information engineering research center, a CCF senior member, a member of IEEE and ACM. His research interests include network and information security, information hiding and steganalysis analysis, image understanding and recognition, multimedia communication technology.



Xing Pengfei Graduated from Henan University of science and Technology, Henan, China, in 2008. He received M.Sc. degrees in Communication systems and communication theory from Lanzhou University of Technology, Lanzhou, China, in 2012. His research interests include audio signal processing and application, multimedia authentication techniques.



Liu Yangwei Received the M.Sc. degrees in Communication systems and communication theory from Lanzhou University of Technology, Lanzhou, China, in 2010. His research interests include audio signal processing and application, multimedia authentication techniques, and steganalysis speech signal processing.