

# Robust Audio Watermarking Algorithm Based on Air Channel Characteristics

Wen Diao, Yuanxin Wu, Weiming Zhang, Bin Liu, Nenghai Yu  
 CAS Key Laboratory of Electromagnetic Space Information  
 University of Science and Technology of China, Hefei, China  
 Email: dw15@mail.ustc.edu.cn, wyx\_ustc@163.com  
 zhangwm@ustc.edu.cn, flowice@ustc.edu.cn ynh@ustc.edu.cn

**Abstract**—This paper proposed an algorithm to improve the robustness of audio watermark for the air channel propagation. Firstly, according to the characteristics of the air channel, noise energy is mainly concentrated in low frequency region. The frequency domain to embed watermark is selected. By modifying the FFT-transformed intermediate frequency coefficients, the watermark is embedded so that the watermark can be robustly extracted while resisting channel interference. The experimental results show that the algorithm has good imperceptibility and can resist 10dB noise interference and resampling attacks. The embedding rate can reach 20bps while maintaining the accuracy rate of more than 92% within 1.5 meters. It can be well adapted to air channel propagation scenario, the performance of the algorithm is greatly improved.

**Index Terms**—audio watermarking, FFT, synchronization signal, air channel

## I. INTRODUCTION

With the rapid development of information technology, the problem, how to protect digital multimedia data against unauthorized copying, is becoming a rising concern. Digital watermarking technology is proposed in this situation where digital works have encountered in the protection of copyright. When digital watermarking technology was proposed, it gained a lot of attention in the industry field immediately. As the main technical means of protecting information security and covert communications, digital watermarking is being widely studied and applied.

The main applications of audio digital watermarking are copyright protection, license control and content authentication. The major performance indicators include: robustness, imperceptibility and payload. Actually, three main indicators are mutually conflicting. Therefore, when designing and developing a new watermarking scheme, it is necessary to seek a balance among the three indicators.

At present, audio piracy products in the field of air channel communication are mainly generated from the recording secretly in public places such as cinemas and concert halls. This makes air channel propagation audio watermarks of great research value. However, in the air propagation, in addition to the complicated DA/AD conversion, the audio watermark will also be affected by the air noise, the volume, distance of the playback source and the orientation of the recording equipment, etc. In addition, audio works may also be subject to re-sampling attacks. Attacks such as format conversion will

cause great technical difficulty for improving the performance of watermarking algorithms.

Previous audio watermarking techniques can be classified into two categories: time domain algorithms [1], [2], [3] and transform domain algorithms [4], [5], [6]. In the time domain algorithm, there are LSB methods [7], [8], echo hiding [9], and so on. The transform domain algorithm hides the watermark by modifying the frequency domain coefficients of the audio signal, such as fast Fourier transform (FFT) [10], discrete cosine transform (DCT) [11] and discrete wavelet transform (DWT) [12]. The most prominent point is that its robustness is significantly enhanced.

Currently, the study of digital audio watermarks for cable channels or for Ethernet transmission is relatively mature. Since there are multiple attacks on air channel transmission process, the research on audio watermark via air-channel is far less than other audio watermark algorithms.

The earliest studies of Steinbach [13], they studied 5 types of audio watermarking technology by using 4 different microphones, but the experimental result is not practical enough. In terms of straight-through cable transmission, Xiang Shijun et al. [14] used a three-stage energy ratio method to embed a string of 32-bit information. Although the performance is relatively good, the capacity is too small, the usability is not good enough and the synchronization technology has higher requirements. Zhang Xia et al. [15] proposed to modify the coefficients in the double DCT domain. Modifying it by selecting a part of the first DCT to perform a second DCT, then embedding the watermark. Although the algorithm can be well concealed, in the noisy environment, bit error rate will be sharp increased.

In the recent work, Andrew Nadeau and Gaurav Sharma [16] proposed an efficient and robust algorithm for resynchronization after analog playback, but it doesn't work in airborne situation. In addition, Michael Arnold et al. [17] used watermark embedding in phase modulation in the WOLA domain, but the embedding rate was low. Qian Wang et al. [18] used the 17KHz and 19KHz frequencies as carriers, which could not be perceived by the human ear, to embed information based on OFDM, and was finally acquired by the smartphone. However, in the actual scene, the high-frequency signal will be erased by filters, which is not suitable for the scene of copyright protection.

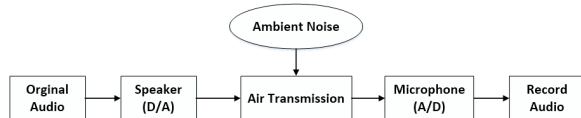


Fig. 1: Audio signal's acoustic transmitting model

It can be seen that how to improve the robustness and concealment effect in public audio broadcasting is still a difficult problem to be solved. In this paper, we propose an audio watermarking algorithm which can effectively resist multiple attacks based on channel characteristics.

The organization of this paper is as follows. The analysis of air channel is presented in Section II. The proposed algorithm is elaborated in Section III. The experimental results are presented in Section IV before concluding this paper in Section V.

## II. THE CHARACTERISTICS OF ACOUSTIC SPEAKER-MICROPHONE CHANNEL

In real life scenarios, most of the air channel propagation comes from the speaker playing audio signals and microphone reception. In order to represent the effect of channel characteristics on audio watermarking better and more specifically. We chose the speaker-microphone channel as a model, and used the experimental framework shown in Fig.1 to explore the characteristics of the channel, illustrating the impact on the watermarked signal.

### A. Impact of DA/AD Conversion Process

Impact of DA/AD conversion process includes noise in the process, linear scaling and waveform distortion on the time axis. The strategy adopted to solve these problems is to add the positioning search synchronization signal to the watermark and embed in the frequency domain.

### B. Impact of Air-transmission Process

The environmental noise in the public places will cause great interference to the audio signal, resulting in a significant increase of error rate during the extracting process. To solve this problem we measure the noise intensity in different environments to characterize these disturbances.

The experimental environment is as follows: We used a Vivo Y19t smartphone to record noise in different environments. The ambient noise energy distribution were measured in square, conference hall and office, as shown in Fig.2. It can be seen that when the frequency is lower than 5Khz, the environmental noise is relatively large. The energy is mainly concentrated in the low-frequency segment. Above 8Khz, the noise is almost negligible.

It gives us an idea of watermark embedding, embedding in the frequency band where the ambient noise is the least disturbing.

At the same time, we play three common types of audio and recorded. As Fig.3 shows, analyzing their spectrum, it

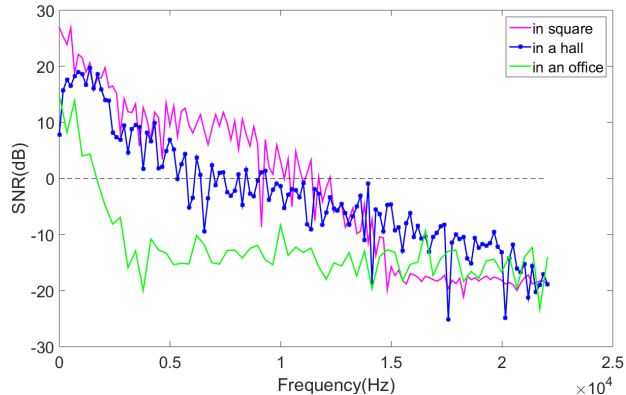


Fig. 2: Spectrum of Ambient Noise

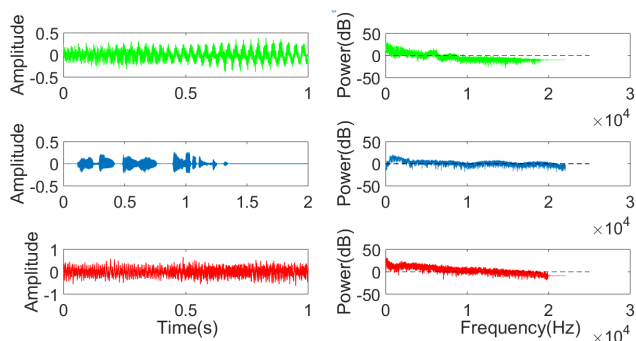


Fig. 3: The waveform and frequency spectrum of pop music, human voice and class music

can be seen that the energy of the audio signal is mainly concentrated in the low and middle frequency. In the process of transmission, the higher frequency is, the more energy decay.

### C. Other impacts

The pre-mute segment which may occur during recording or in original audio.

## III. PROPOSED AUDIO WATERMARKING ALGORITHM

Because of the linear stretching, waveform distortion and noise interference and the presence of silent segments, we adopt corresponding countermeasures in the embedding and extraction algorithms respectively:

- 1) Audio preprocessing: As shown in Fig.4, using double-threshold method [19] to find a suitable processing start point.
- 2) Adding synchronization code in the watermark embedding process, which is used to resist the linear expansion in the channel propagation process, and can eliminate the linear expansion and contraction in the extraction process.
- 3) To resist the waveform distortion and noise interference in the propagation process, we perform FFT transform on the signal and modify the coefficients in the mid-band of the FFT domain to achieve the purpose of embedding the watermark.

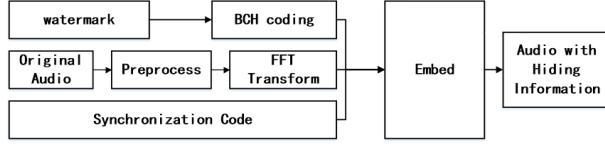


Fig. 4: The framework of embedding algorithm

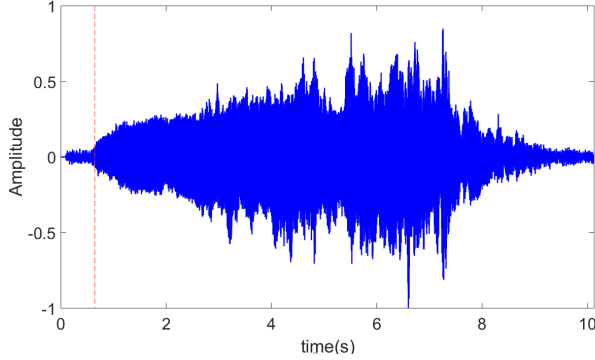


Fig. 5: Using endpoint detection to find the starting point

#### A. Audio Preprocessing

In the audio signal, there will be a period of silence or very small amplitude. At the same time, pre-mute may still occur during recording, synchronization marks and watermark embedding effects are generally performed badly in this area. For example, the situation shown in the Fig.5.

There is a situation frequently encountered in this area, such that the identification bit is not found and the decoding error is found. In order to reduce the occurrence of this kind of situation, it is necessary to pretreat the audio signal firstly. Using the endpoint detection technology to detect the right location for follow-up operations.

To improve the quality of the embedding, double threshold detection is used to determine the starting point  $Start_p$ .

#### B. Watermark embedding algorithm

According to previous experiments and analysis, it can be judged that modifying the frequency band in FFT can effectively resist noise attacks and other universal attacks.

Original audio signal is divided into segments. Every segment includes a synchronization frame and several watermark frames. First, the watermark message is BCH-encoded [20]. The encoded watermark message is marked as  $w(i)$ . Then we perform FFT transformation on each watermark frame, and modify the coefficients to embed watermark message.

The FFT transform equation is shown as (1), the Inverse FFT transform is shown as (2).

$$f(k) = \sum_{n=0}^{N-1} x(n) \cdot e^{-\frac{2\pi i}{N} \cdot kn}, \quad k = 0, \dots, N-1 \quad (1)$$

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} f(k) \cdot e^{-\frac{2\pi i}{N} \cdot kn}, \quad n = 0, \dots, N-1 \quad (2)$$

In this paper, We define the frequencies from 8K to 10K as intermediate frequency. Selecting a frequency in intermediate frequency randomly as the centre point, marked as  $f_1$ . We take  $L$  points around  $f_1$ , marked as  $f(0), f(1), f(2) \dots f(L-1)$  and divide it into two segments. Each segment is  $L/2$  long. The energy of the two-stage FFT coefficients is  $E_1, E_2$  is defined as follows.

$$E_1 = \sum_{i=0}^{L/2-1} |f(i)| \quad (3)$$

$$E_2 = \sum_{i=L/2}^{L-1} |f(i)| \quad (4)$$

Simultaneously, setting the embedded strength  $S$ ,  $S = \alpha \cdot \max(E_1, E_2)$ , where  $\alpha$  is the strength factor. To resist interference during the air transmission, the  $\alpha$  value should be as large as possible under the constraint of imperceptibility. The parameter  $\alpha$  is assigned as a predefined value at the beginning, and then automatically adjusted until the objective quality grade (ODG) value of watermarked audio is satisfied. In the proposed strategy, one watermark bit  $w(i)$  can be embedded by modifying the relationships among  $E_1, E_2$  and  $S$ , as shown in (5):

$$\begin{cases} E_1 - E_2 \geq S, & \text{if } w(i) = 1 \\ E_2 - E_1 > S, & \text{if } w(i) = 0 \end{cases} \quad (5)$$

Basic steps involved in the watermarking embedding are given as follows. We give an example as watermark  $w(i)=1$  to illustrate.

**Step 1:** If  $E_1 - E_2 \geq S$ , coefficients do not adjust,  $f(i)' = f(i)$ . If not, do step 2.

**Step 2:** To increase  $E_1$  and decrease  $E_2$ , the specific measure is

$$f(i)' = f(i) + \delta, \quad i = 0, 1, 2, \dots, L/2-1 \quad (6)$$

and

$$f(i)' = f(i) - \delta, \quad i = L/2, L/2+1, \dots, L \quad (7)$$

$\delta$  is a small non-negative value,  $f(i)'$  is a modified value.

**Step 3:** Case 1,  $E_1 - E_2 \geq S$ , stop. Case 2,  $E_1 - E_2 < S$ , return step 2.

Original FFT coefficients are shown in Fig.6, watermarking result is shown in Fig.7.

Finally, because the coefficients of the FFT transform are conjugate symmetric, the same operation should be performed at the symmetry point of the FFT transform domain, then IFFT is performed to obtain the time domain information of the embedded watermark information.

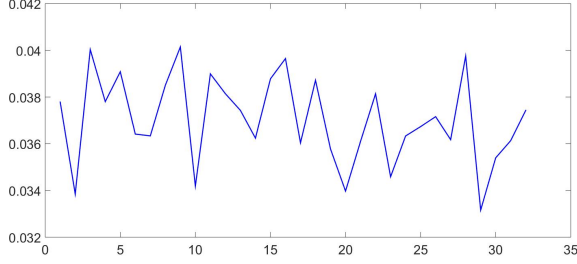


Fig. 6: Coefficients of watermark frame in FFT domain

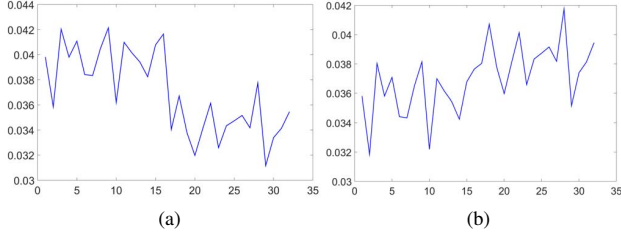


Fig. 7: Watermarking '1'(a) and Watermarking '0'(b) coefficients in FFT domain

### C. Synchronization Code Embedding

According to our previous test results, embedding synchronization signals in the midband of the audio can be effective against air interference and background noise in the surrounding environment. We use a hidden information structure of the synchronization signal and the watermark signal loop body, its structure is shown in Fig.8.

The length of the synchronization signal is  $N_1$ , and the length of the watermark signal is  $N_2=15 \times N_1$ . We perform a FFT on the synchronization signal frame of  $N_1$  to calculate the absolute value as its energy. Select another frequency in intermediate frequency  $f_2$  different from  $f_1$ . Mark the maximum energy as  $Max\_f$  and take the 8 points  $c(i)$ ,  $0 \leq i \leq 7$ , on both sides of  $f_2$  position where the maximum value marked as  $Max\_c$ , changing it as (8) shown:

$$c(i)' = \beta \cdot \frac{c(i)}{Max\_c} \cdot Max\_f \quad (8)$$

Among them,  $\beta$  is the modification factor, which can be modified according to the audio quality requirements. Here we set it to be 0.5. Same as the embedding process, we take 8 points in the symmetry positions and do the same embedding operation. The original signal's sepectrogram is shown in Fig.9. The embedded spectrum is shown as Fig.10. Finally, the

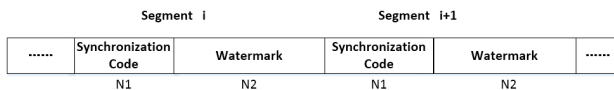


Fig. 8: Construction of embedding information

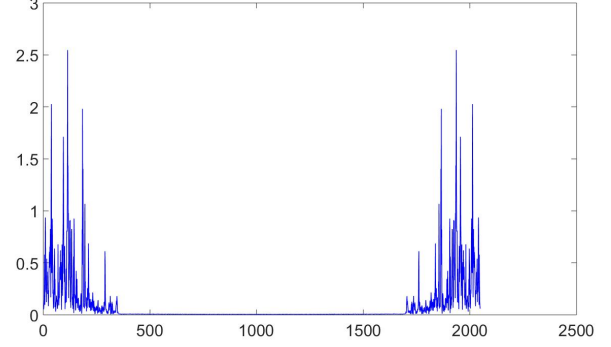


Fig. 9: FFT Spectrogram of Original signal

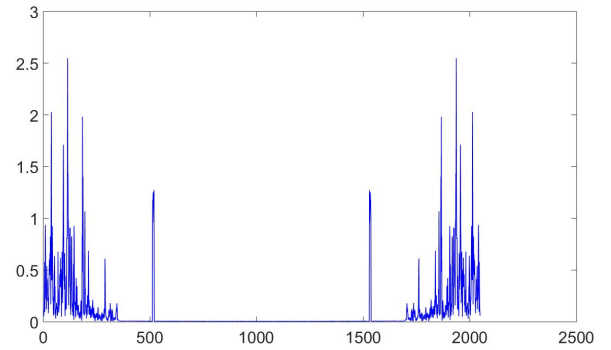


Fig. 10: FFT Spectrogram of embedded syn signal

frequency domain signal is IFFT transformed to obtain the embedded time domain information.

### D. Watermarking extracting algorithm

The watermark extraction algorithm is as Fig.11 shown. Firstly, the audio signal is preprocessed and then synchronously decoded. The location of the watermark message is determined by the location of the synchronization code. Then the FFT transform is performed, the FFT coefficients are compared to extract the message, and finally the BCH decoding is performed to obtain the watermark message. The specific process is as Fig.12.

1) *Synchronization decoding*: First, we perform synchronous decoding. The sliding window structure shown in Fig.12 is used for decoding. Setting the threshold  $t\_value$ . We put  $N_1$  points as a frame, then do FFT tranform and take the absolute value. The maximum value is recorded as  $Max\_f_0$ .

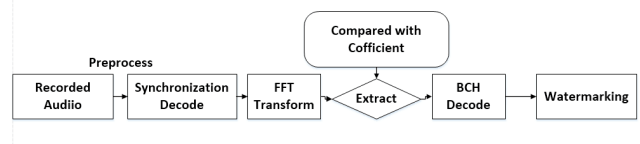


Fig. 11: Framework of extraction algorithm

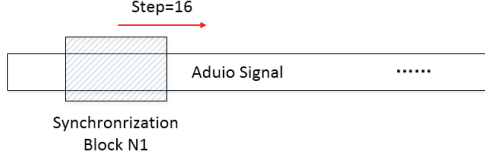


Fig. 12: Framework of resynchronization algorithm

We take the 8 points on both sides of  $f_2$ , mark their values as  $c(i)$ .  $i=0,1,\dots,7$ . Calculating the sum of 8 points energy mark as  $sum\_x$ , if

$$sum\_x / Max\_f_0 \geq t\_value \quad (9)$$

one synchronization information in the signal is found. Recording its corresponding position. Because we embedded synchronization signal in the way of circulation, we skip the following section of watermarking information, and synchronize again. If

$$sum\_x / Max\_f_0 < t\_value \quad (10)$$

the window is slid for 16 samples and synchronization is performed again until the end of the signal.

2) *Watermark detection*: After finding the synchronization signal in 1) and determining the embedding position, it is processed according to the method of embedding, each watermark frame is subjected to FFT transformation. The  $L$  points around the intermediate frequency  $f_1$  are taken and divided into two segments to compare the energy levels.

$$E'_1 = \sum_{i=0}^{L/2-1} |f(i)| \quad (11)$$

$$E'_2 = \sum_{i=L/2}^{L-1} |f(i)| \quad (12)$$

and

$$\begin{cases} E'_1 \geq E'_2, w'(i) = 1; \\ E'_1 < E'_2, w'(i) = 0; \end{cases} \quad (13)$$

In this time, we get a hidden sequence  $w'(i)$ . After the BCH error correction in the decoding of sequence, we get the original watermark.

### E. Performance evaluation

In this part, we will discuss the algorithm's embedding rate and signal-to-noise ratio(SNR), and demonstrate that the algorithm is feasible.

#### 1) Evaluation of Embedding capacity:

$$B = \frac{fs}{N_1 + N_2} \cdot \frac{N_2}{N_1} \quad (bps) \quad (14)$$

In this equation,  $fs$  represents signal sampling rates.  $N_1 = 2048$  represents the length of synchronization frame.  $N_2$  represents the length of embedding frame. Here, sampling rate is 44.1Khz and  $N_2$  equal to 15 times of  $N_1$ , the embedding capacity is 20.18bps.

#### 2) Evaluation of SNR:

$$\begin{aligned} SNR &= -10lg \left( \frac{\|F-F'\|^2}{\|F\|^2} \right) \\ &= -10lg \left( \frac{\sum_{i=0}^{N-1} (f(i)^2) - \sum_{i=0}^{N-1} (f(i)')^2}{\sum_{i=0}^{N-1} (f(i)^2)} \right) \end{aligned} \quad (15)$$

In (15),  $F$  is original signal, and  $F'$  is watermarked audio signal.  $f(i)$  is original FFT coefficient and  $f(i)'$  is modified coefficient in FFT domain. Experiment results indicate the average SNR is greater than 20dB, which is satisfied with the Sound of International Union's requirement.

## IV. TEST AND ANALYSIS PERFORMANCE

In order to verify the robustness of the algorithm and the applicability in the real scene, we did a lot of experiments. The results and details are as follows.

### A. The impact of distance

In the air channel, distance plays an important role as a indicator. Ambient noise is also a very important indicator. In order to test the effect of distance and noise to the robustness of algorithm, we have selected four typical audio frequencies. A bit error rate (BER) test was performed at distances of 0.2m, 0.5m, 1m, and 1.5m.

As Fig.13 shown, the experiment now indicates that the algorithm is robust enough and the correct rate is effectively improved at relatively close distances.

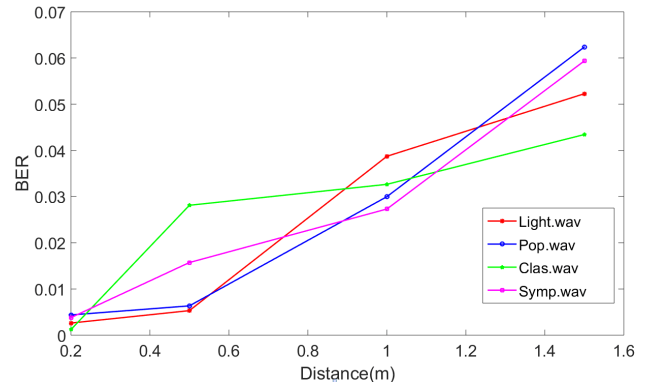


Fig. 13: The impact of distance between speaker and microphone

### B. The impact of noise

Propagating in the air is unavoidably interfered by noise. In order to detect the anti-noise ability of the algorithm, we have done a bit error rate test at different SNRs and compared it with algorithms [14] and [15]. The specific results are shown in Fig.14.

The test result shows that the algorithm has strong robustness and can resist the interference of different SNR noises, even 10dB. Compared with the previous algorithms, the performance of the algorithm is greatly improved, especially at low SNR.

