

# DIGITAL WATERMARKING OF AUDIO SIGNALS USING A PSYCHOACOUSTIC AUDITORY MODEL AND SPREAD SPECTRUM THEORY

**RICARDO A. GARCIA\***, *AES Member*

*School of Music Engineering Technology, University of Miami  
Coral Gables, FL 33146, USA*

*A new algorithm for embedding a digital watermark into an audio signal is proposed. It uses spread spectrum theory to generate a watermark resistant to different removal attempts and a psychoacoustic auditory model to shape and embed the watermark into the audio signal while retaining the signal's perceptual quality. Recovery is performed without knowledge of the original audio signal. A software system is implemented and tested for perceptual transparency and data-recovery performance.*

## 0 INTRODUCTION

Every day the amount of recorded audio data and the possibilities to distribute it (i.e. by the Internet, CD recorders, etc) are growing. These factors can lead to an increase in the illicit recording, copying and distributing of audio material without respect to the copyright or intellectual property of the legal owners. Another concern is the tracking of audio material over broadcast media without the use of human listeners or complicated audio recognition devices. Audio watermarking techniques promise a solution to some of these problems.

The concept of watermarking has been used for years in the fields of still and moving images. The basic idea of a watermark is to include a special “code” or information within the transmitted signal. This code should be transparent to the user (non-perceptible) and resistant against removal attacks of various types.

In audio signals, the desired characteristics can be translated into:

- Not perceptible (the audio information should appear “the same” to the average listener before and after the code is embedded).
- Resistant to degradation because of analog channel transmission. (i.e. TV, radio and tape recording).
- Resistant to degradation because of uncompressed-digital media. (i.e. CD, DAT and wav files).
- Resistant to removal through the use of sub-band coders or psychoacoustic models. (i.e. MPEG, Atrac, etc).

The proposed algorithm generates a digital watermark (i.e. a bit stream) that is spectrally shaped and embedded into an audio signal. Spread spectrum theory is used in the generation of

---

\* Currently with the program in Media Arts & Sciences, Machine Listening Group, Massachusetts Institute of Technology (MIT), Cambridge, MA 02139 – 4307 USA.

the watermark. The strength of coded direct-sequence/binary-phase-shift keying (DS/BPSK) is used to create a robust watermark. The concepts are adapted to better deal with audio signals in a restricted audio bandwidth. A psychoacoustic auditory model is applied to shape and embed the watermark into the audio signal while retaining its perceptual quality for the average listener.

A complete psychoacoustic auditory model algorithm is explained in detail. This information is useful for other applications involving auditory models. The spread spectrum encoding and decoding processes are then presented. The algorithm performs an analysis of the incoming signal and searches the frequency domain for “holes” in the spectrum where the spread spectrum data can be placed without being perceived by the listener. The psychoacoustic auditory model is used to find these frequency “holes.”

After transmission, the receiver recovers the embedded spread spectrum information and decodes it in order to reconstruct the original bit stream (watermark). There is no need for the receiver to have access to the original audio signal.

The algorithm was implemented in a software system to create an encoder and decoder, and its performance was evaluated for diverse channels and audio signals. The survival of the watermark (number of correct bytes/second) was analyzed for different configurations of the encoding system. Each one of these configurations was tested for transparency using an ABX listening test and for different channels (i.e. AM Radio, FM stereo radio, Mini Disc, MPEG layer 3, D/A – A/D conversion, etc).

## **1 PSYCHOACOUSTIC AUDITORY MODEL**

An auditory model is an algorithm that tries to imitate the human hearing mechanism. It uses knowledge from several areas such biophysics and psychoacoustics.

From the many phenomena that occur in the hearing process, the one that is the most important for this model is “simultaneous frequency masking.” The auditory model processes the audio information to produce information about the final masking threshold. The final masking threshold information is used to shape the generated audio watermark. This shaped watermark is ideally imperceptible for the average listener. To overcome the potential problem of the audio signal being too long to be processed all at the same time, and also extract quasi-periodic sections of the waveform, the signal is segmented in short overlapping segments, processed and added back together. Each one of these segments is called a “frame.”

The steps needed to form a psychoacoustic auditory model are condensed in Figure 1. The first step is to translate the actual audio frame signal into the frequency domain using the Fast Fourier Transform. In the frequency domain the power spectrum, energy per critical band and the spread energy per critical band are calculated to estimate the masking threshold. This masking threshold is used to shape the “noise or watermark” signal to be imperceptible (below the threshold). Finally frequency domain output is translated into the time domain and the next frame is processed.

### **1.1 Short Time Fourier Transform (STFT)**

The cochlea can be considered as a mechanical to electrical transducer, and its function is to make a time to frequency transformation of the audio signal. To be more specific, the audio information, in time, is translated in first instance into a frequency-spatial representation inside the basilar membrane. This spatial representation is perceived by the nervous system and translated into a frequency-electrical representation.

This phenomenon is modeled using the short time Fourier Transform (STFT). The STFT uses successive, overlapped windows from the time domain input signal.

## 1.2 Simultaneous Frequency Masking and Bark Scale

Simultaneous masking of sound occurs when two sounds are played at the same time and one of them is masked or “hidden” because of the other. The formal definition says that masking occurs when a test tone or “maskee” (usually a sinusoidal tone) is barely audible in the presence of a second tone or “masker.” The difference in sound pressure level between the masker and maskee is called the “masking level.” [ 1 ]

It is easier to measure the masking level for narrow band noise maskers (with a defined center frequency) and sinusoidal tone maskees. Figure 2 (a) and (b) display some curves that show the masking threshold for different narrow band noise maskers centered at 70, 250, 1000 and 4000 Hz. The level of all the maskers is 60 dB. The broken line represents the “threshold in quiet.” Average listeners will not hear any sound below this threshold. Figure 2 (a) uses a linear and (b) uses a logarithmic frequency scales.

The shape of all the masking curves is very different across the frequency range in both graphs. There are some similarities in the shape of the curves below 500 Hz in the linear frequency scale (a), and some similarities above 500 Hz in the logarithmic frequency scale (b). A more useful scale has been introduced that is known as “critical band rate” or “Bark scale.” The concept of the Bark scale is based on the well-researched assumption [ 1 ] that the basilar membrane in the hearing mechanism analyzes the incoming sound through a spatial-spectral analysis. This is done in small sectors or regions of the basilar membrane that are called “critical bands.” If all the critical bands are added together in a way that the upper limit of one is the lower limit of the next one, the critical band rate scale is obtained. Also a new unit has been introduced, the “Bark” that is by definition one critical band wide.

Figure 3 shows the same masking curves from Figure 2 in a Bark scale. Notice that the shape of the masking curves is almost identical across the frequency range. Various approximations may be used to translate frequency into a Bark scale [ 2 ]:

$$z = 13 \tan^{-1}\left(\frac{0.76 * f}{1000}\right) + 3.5 \tan^{-1}\left(\left(\frac{f}{7500}\right)^2\right) \quad (1)$$

and [ 3 ]:

$$z = \frac{26.81 * f}{1960 + f} - 0.53 \quad (2)$$

where  $f$  is the frequency in Hertz and  $z$  is the mapped frequency in Barks.

Eq. ( 1 ) is more accurate, but the Eq. ( 2 ) is easier to compute. Figure 3 shows the excitation level of several narrow band noises with diverse center frequencies in a Bark scale.

## 1.3 Power Spectra

The first step in the frequency domain (linear, logarithmic or bark scales) is to calculate the power spectra of the incoming signal. This is calculated with:

$$\begin{aligned} Sp(j\omega) &= \text{Re}\{Sw(j\omega)\}^2 + \text{Im}\{Sw(j\omega)\}^2 \\ &= |Sw(j\omega)|^2 \end{aligned} \quad (3)$$

The energy per critical band,  $Spz(z)$ , is defined as:

$$Spz(z) = \sum_{w=LBZ}^{HBZ} Sp(jw) \quad (4)$$

Where:  $z = 1, 2, \dots$ , total of critical bands  $Zt$ ;  $LBZ$  and  $HBZ$  the lower and higher frequencies in the critical band  $z$ .

The power spectrum  $Sp(jw)$  and the energy per critical band  $Spz(z)$  are the base of the analysis in the frequency domain. They will be used to compute the spread masking threshold.

## 1.4 Basilar Membrane Spreading Function

A model that approximates the basilar membrane spreading function, without taking in account the change in the upper slope is defined [ 3 ]:

$$B(z) = 15.91 + 7.5(z + 0.474) - 17.5\sqrt{1 + (z + 0.474)^2} \quad (5)$$

where  $z$  is the normalized Bark scale. Figure 4 shows  $B(z)$ .

The auditory model uses the information about the energy in each critical band given by Eq. ( 4 ) and uses Eq. ( 5 ) to calculate the spread masking across critical bands  $Sm(z)$ . This is done using:

$$Sm(z) = Spz(z) * B(z) \quad (6)$$

This operation is a convolution between the basilar membrane spreading function and the total energy per critical band. A true spreading calculation should include all the components in each critical band, but for the purposes of this algorithm, the use of the energy per critical band  $Spz(z)$  is a close approximation.  $Sm(z)$  can be interpreted as the energy per critical band after taking in account the masking occasioned by neighboring bands.

## 1.5 Masking Threshold Estimate

### 1.5.1 Masking Index

There are two different indexes used to model masking. The first one is used when a tone is masking noise (masker = tone, maskee = noise), and it is defined to be  $14.5 + z$  dB below the spread masking across critical bands  $Sm(z)$ . In this case  $z$  is the center frequency of the masker tone using a bark scale. The second index is used when noise is masking a tone (masker = noise, maskee=tone), and is defined to be 5.5 dB below  $Sm(z)$ , regardless of the center frequency [ 4 ].

### 1.5.2 Spectral Flatness Measure (SFM) and Tonality Factor $\alpha$

The spectral flatness measure (SFM) is used to determine if the actual frame is noise-like or tone-like and then to select the appropriate masking index. The SFM is defined as the ratio of the geometric to the arithmetic mean of  $Spz(z)$ , expressed in dB as:

$$SFM_{dB} = 10 \log_{10} \left\{ \frac{\prod_{z=1}^{Zt} Spz(z)}{\frac{1}{Zt} \sum_{z=1}^{Zt} Spz(z)} \right\}^{\frac{1}{Zt}} \quad (7)$$

with  $Zt$  = total number of critical bands on the signal

The value of the SFM is used to generate the “tonality factor” that will help to select the right masking index for the actual frame. The tonality factor is defined in [ 3 ], [ 4 ] as the minimum of the ratio of the calculated SFM over a SMF maxima and 1:

$$\mathbf{a} = \min\left(\frac{SFM_{dB}}{SFM_{dB_{max}}}, 1\right) \quad (8)$$

with  $SFM_{dB_{max}} = -60dB$ .

Therefore, if the analyzed frame is tone-like, the tonality factor  $\alpha$  will be close to 1, and if the frame is noise-like,  $\alpha$  will be close to 0. The tonality factor  $\alpha$  is used to calculate the masking energy offset  $O(z)$ , is defined as [ 3 ], [ 4 ]:

$$O(z) = \mathbf{a}(14.5 + z) + (1 - \mathbf{a})5.5 \quad (9)$$

The offset  $O(z)$  is subtracted from the spread masking threshold to estimate the raw masking threshold  $Traw(z)$ .

$$Traw(z) = 10^{\left(\log_{10}(Sm(z)) - \frac{O(z)}{10}\right)} \quad (10)$$

### 1.5.3 Threshold Normalization

The use of the spreading function  $B(z)$  increases the energy level in each one of the critical bands of the spectrum  $Sm(z)$ . This effect has to be undone using a normalization technique, to return  $Traw(z)$  to the desired level. The energy per critical band calculated with Eq. ( 4 ) is also affected by the number of components in each critical band. Higher bands have more components than lower bands, affecting the energy levels by a different amount. The normalization used [ 4 ] simply divides each one of the components of  $Traw(z)$  by the number of points in the respective band  $P_z$ .

$$Tnorm(z) = \frac{Traw(z)}{P_z} \quad (11)$$

Where:  $P_z$  = number of points in each band  $z$   
 $z = 1, 2, \dots, Z_t$

### 1.5.4 Final Masking Threshold

After normalization, the last step is to take in to account the absolute auditory threshold or “hearing threshold.” The hearing threshold varies across the frequency range as stated in Zwicker and Zwicker [ 1 ]. In the proposed auditory model the hearing threshold will be simplified to use the worst case threshold (the lowest). That is defined as a sinusoidal tone of 4000 Hz with one bit of dynamic range [ 4 ]. These values are chosen based on the data from experimental research that shows that the most sensitive range of the human ear is in the range of 2500 to 4500 Hz [ 1 ]. For a frequency of 4000 Hz, the measured sound intensity is  $10^{-12}$  Watt/m<sup>2</sup>, that equals a loudness of 0 phons at that frequency [ 12 ]. The chosen amplitude (one bit) is the smallest possible amplitude value in a digital sound format. The hearing threshold is then calculated with [ 4 ]:

$$TH = \max(|Pp(j\omega)|) \quad (12)$$

where:  $Pp(j\omega)$  = power spectrum of the probe signal  $p(t)$   
 $p(t) = \sin(2\pi 4000t)$

The final threshold  $T(z)$  is:

$$T(z) = \max(T_{norm}(z), TH) \quad (13)$$

## 1.6 Noise Shaping Using the Masking Threshold

The objective of the auditory model is to find a usable masking threshold. The final masking threshold is always compared with the values of the power spectrum of the signal  $Sp(j\omega)$ . This can be interpreted as “below this threshold, the information is not relevant for human hearing.” This means that if the frequency components that fall below the masking threshold are removed; the average listener will notice no difference between the original sound signal and the altered version.

Another very important consequence of this is that if these components are not just discarded but replaced with new components they will be, as before, inaudible for the listener. This assumes that the new components do not change the average energy considerably in their critical band. Let the frame with the new components be called  $N(j\omega)$ . The objective is to use the final masking threshold to select which components from  $Sp(j\omega)$  can be replaced with components from  $N(j\omega)$ . The components of  $N(j\omega)$  are shaped to stay below the final masking threshold. The final signal, that includes components from  $Sw(j\omega)$  and  $N(j\omega)$ , ideally retains the perceptual quality of the original signal for the average listener.

The following steps are used to remove the components from  $Sw(j\omega)$ , shape the vector  $N(j\omega)$  and mix them:

Calculate the “new” version of the sound signal (after removing some components):

$$S_{wnew_i}(j\omega) = \begin{cases} Sw_i(j\omega) & Sp_i(j\omega) \geq T(z) \\ 0 & Sp_i(j\omega) < T(z) \end{cases} \quad (14)$$

$i = 1, 2, \dots$  number of components  
 $z, \omega$  according to component  $i$

Remove the unneeded components in the  $N(j\omega)$  vector:

$$N_{new_i}(j\omega) = \begin{cases} 0 & Sp_i(j\omega) \geq T(z) \\ N_i(j\omega) & Sp_i(j\omega) < T(z) \end{cases} \quad (15)$$

$i = 1, 2, \dots$  number of components  
 $z, \omega$  according to component  $i$

Calculate the power spectrum of  $N_{new}(j\omega)$ :

$$N_{newp}(j\omega) = |N_{new}(j\omega)|^2 \quad (16)$$

and then, the energy per critical band:

$$N_{newpz}(z) = \sum_{\omega=LBZ}^{HBZ} N_{newp}(j\omega) \quad (17)$$

Where:  $z = 1, 2, \dots$ , total of critical bands  $Z_t$ ; LBZ and HBZ the lower and higher frequencies in the critical band  $z$ .

The shaping is done applying a factor  $F_z$  to each critical band. These factors are given by:

$$F_z = A \frac{\sqrt{T(j\omega)}}{\max(|N_{new}(j\omega)|)} \quad (18)$$

$$z = 1, 2, \dots, Z_t$$

$\omega = LBZ$  to  $HBZ$  for each band  $z$

The coefficient  $A$  is used as the “gain of the noise signal”. Varies from 0 to 1 and weights the embedded noise below the threshold of masking. The factors  $F_z$  are applied using:

$$N_{final}(j\omega) = N_{new}(j\omega)F_z \quad (19)$$

$$z = 1, 2, \dots, Z_t$$

$\omega = LBZ$  to  $HBZ$  for each band  $z$

The final step is to mix both spectrums, the altered  $S_{new}(j\omega)$  and the shaped  $N_{final}(j\omega)$  to form the composite signal  $OUT(j\omega)$ :

$$OUT(j\omega) = S_{new}(j\omega) + N_{final}(j\omega) \quad (20)$$

## 2 SPREAD SPECTRUM

One of the requirements of a watermarking algorithm is that the watermark should resist multiple types of removal attacks. A removal attack is considered as anything that can degrade or destroy the embedded watermark. Another factor to be considered is that the masking threshold of the actual audio signal determines the embedding of the watermark, because the watermark is embedded in the “spare components” found using the psychoacoustic auditory model. From this point of view, the watermark has to be the least intrusive to the audio signal, and therefore, the actual audio data can be seen as the main obstacle for a good watermarking algorithm. This is because the audio will use all the needed bandwidth and the watermark will use what is left after the auditory model analysis.

The desired watermarking technique should be resistant to degradation because of:

- The used transmission channel: analog or digital.
- High-level wide-band noise (in this case, the “noise” is the actual audio signal). This is often related as “low signal to noise ratio”.
- The use of psychoacoustic algorithms on the final watermarked audio.

A communication theory technique that meets the requirements is the “spread spectrum technique”, as described thoroughly in Simon et al. [ 5 ] and Pickholtz et al.

[ 6 ]. “*Spread spectrum is a means of transmission in which the signal occupies a bandwidth in excess of the minimum necessary to send the information; the band spread is accomplished by means of a code which is independent of the data, and a synchronized reception with the code at the receiver is used for despreading and subsequent data recovery.*” [ 6 ]

In the following analysis, the process of generating a watermark that will be embedded in an audio signal is expressed in spread spectrum terminology. The original audio signal will be called “noise” and the bit stream that conforms the watermark sequence will be the data signal. The watermark sequence is transformed in a watermark audio signal and then the audio signal (noise) is added to it. This process of adding noise to a channel or signal is called “jamming.” The objective of a jammer in a communication system is to degrade the performance of the transmission, exploiting knowledge of the communication system. In the watermark algorithm

the audio signal (i.e. music) is considered the jammer, and it has much more power than the transmitted bit stream (watermark).

## 2.1 Basic Concepts

The primary challenge that a receiver must overcome is intentional jamming, especially if the jammer has much more power than the transmitted signal. Classical communications theoretical investigations about additive white Gaussian noise help to analyze the problem. White Gaussian noise is a signal which has infinite power spread uniformly over all frequencies; but even under these circumstances communication can be achieved due to the fact that on each of the “signal coordinates” the power of the noise component is limited (not infinite). Therefore, if the noise component in the signal coordinates is not too large, communication can be made. This is usually applied in a typical narrow-band signal, where just the noise components in the signal bandwidth are taken into account as possible factors that can do harm to the communication. With this knowledge, the best strategy to combat intentional jamming is to select signal coordinates where the jammer to signal ratio is the smallest possible.

Assume a communication link with many signal coordinates available to choose from, and only a small subset of these is used at any time. If the jammer can not determine which subset is being used, it is forced to jam all the coordinates and therefore, all its power will be distributed among all the coordinates, with little power in each of them. If the jammer chooses to jam only some of the coordinates, the power over each of them is larger, but the jammer lacks the knowledge of which coordinates to jam. The protection against the jammer is enhanced, as more signal coordinates are available to choose from.

Having a signal of bandwidth  $W$  and duration  $T$ , the number of coordinates available is given by:

$$N \cong \begin{cases} 2WT & \text{coherent signals} \\ WT & \text{non-coherent signals} \end{cases} \quad (21)$$

$T$  is the time used to send a standard symbol. To make  $N$  larger when  $T$  is fixed, two techniques can be applied:

- Direct sequence spreading (DS): this is the selected approach in this algorithm.
- Frequency hopping (FH)

The signals created with these techniques are called “spread spectrum signals.”

### 2.1.1 Models and Fundamental Parameters

The basic system is shown in Figure 5, with the following parameters:

$W_{ss}$  = Total spread spectrum signal bandwidth available

$R_b$  = Data rate ( bits / second )

$S$  = Signal power (at the input of the receiver)

$J$  = Jammer power (at the input of the receiver)

$W_{ss}$  is defined as the total available spread spectrum bandwidth that could be used by the transmitter, but it is not guaranteed that it will be used during the actual transmission. Neither is it guaranteed that the spectrum will be continuous.  $R_b$  is the uncoded bit data rate used during transmission. The signal and the jammer powers  $S$  and  $J$  are the averaged power at the receiver. This does not change even if the jammer and/or the signal are pulsating.

### 2.1.2 Jammer Waveforms

The number of possible jammer waveforms that a jammer can apply to a communication system is infinite. The principal types include:

- *Broadband Noise Jammer*: Spreads Gaussian noise of a total power  $J$  evenly over the total frequency range of the spread bandwidth  $W_{ss}$ .
- *Partial Band Noise Jammer*: Spreads noise of total power  $J$  evenly over a frequency range of bandwidth  $W_J$ , which is contained in the total spread bandwidth  $W_{ss}$ .  $r$  is the fraction of the total spread spectrum bandwidth that is being jammed.
- *Pulse Jammer*: Transmits the jammer waveform during a fraction  $r$  of the time, the average power is  $J$ , but the peak power during transmission is higher.

## 2.2 Coherent Direct-Sequence Systems

Coherent direct-sequence systems use a pseudorandom sequence and a modulator signal to modulate and transmit the data bit stream. The main difference between the uncoded and coded versions is that the coded version uses redundancy and “scrambles” the data bit stream before the modulation is done and reverses the process at the reception. The watermarking algorithm uses the coded scheme, but the uncoded is studied because is easier to understand and is the foundation of the coded scheme.

### 2.2.1 Uncoded Direct-Sequence Spread Binary Phase-Shift-Keying

Uncoded Direct-Sequence Spread Binary Phase-Shift-Keying is known as Uncoded DS/BPSK. It may be explained with a simple example. BPSK signals are often expressed as:

$$s(t) = \sqrt{2S} \sin \left[ \omega_0 t + \frac{d_n \pi}{2} \right] \quad (22)$$

$$nT_b \leq t < (n+1)T_b, \quad n = \text{integer}$$

where

$$T_b \text{ is the data bit time } \left( \frac{1}{R_b} \right)$$

$\{d_n\}$  is the sequence of data bits, with the possible values of 1 or  $-1$ ; and equal probability of occurrence.

Eq. ( 22 ) can be expressed as:

$$s(t) = d_n \sqrt{2S} \cos(\omega_0 t) \quad (23)$$

$$nT_b \leq t < (n+1)T_b, \quad n = \text{integer}$$

BPSK can be seen as phase modulation in Eq.( 22 ) or amplitude modulation in Eq. ( 23 ). The spectrum of a BPSK signal is usually of the form shown in Figure 6. This is a  $(\sin^2 x)/x^2$  function, and the first null bandwidth is  $1/T_b$ . This shows the minimum bandwidth needed to transmit the signal  $s(t)$  and to recover it at the receiver.

Spread spectrum theory requires the signal to be spread over a larger spectrum than the minimum needed for transmission. The spreading of the direct sequence is done using a pseudorandom (PN) binary sequence  $\{c\}$ . The values of this sequence are 1 or  $-1$  and its speed is  $N$  times faster than the  $\{d\}$  data rate. The time,  $T_c$ , of each bit on a PN sequence is known as a “chip” and is given by:

$$T_c = \frac{T_b}{N} \quad (24)$$

The direct sequence spread spectrum signal has the form:

$$\begin{aligned} x(t) &= \sqrt{2S} \sin[\mathbf{w}_0 t - d_n c_{nN+k} \mathbf{p} / 2] \\ &= d_n c_{nN+k} \sqrt{2S} \cos(\mathbf{w}_0 t) \\ nT_b + kT_c &\leq t < nT_b + (k+1)T_c \\ k &= 0,1,2 \dots N-1 \\ n &= \text{integer} \end{aligned} \quad (25)$$

The signal is very similar to the common BPSK, except that the bit rate is  $N$  times faster and the power spectrum is  $N$  times wider, as shown in Figure 7. The processing gain is given by:

$$PG = \frac{W_{SS}}{R_b} = N \quad (26)$$

$W_{SS}$  is the direct sequence spread spectrum bandwidth  $\frac{1}{T_c} = N \frac{1}{T_b}$ .

If the data function is defined as:

$$\begin{aligned} d(t) &= d_n, \quad nT_b \leq t < (n+1)T_b \\ n &= \text{integer} \end{aligned} \quad (27)$$

and the PN sequence is:

$$\begin{aligned} c(t) &= c_k, \quad kT_c \leq t < (k+1)T_c \\ k &= \text{integer} \end{aligned} \quad (28)$$

Eq. ( 25 ) can be expanded as:

$$\begin{aligned} x(t) &= \sqrt{2S} \sin[\mathbf{w}_0 t + c(t)d(t)\mathbf{p} / 2] \\ &= c(t)d(t)\sqrt{2S} \cos(\mathbf{w}_0 t) \end{aligned} \quad (29)$$

Figure 8 shows the block diagram for the normal DS/BPSK modulation; and Figure 9 shows an equivalent model used in the next step of the analysis. Figure 11 shows the signals  $d(t)$  and  $c(t)$  and Figure 12 shows  $c(t)d(t)$  with  $N=6$ . From Figure 9, the equivalent form of  $x(t)$  is given by:

$$x(t) = c(t)s(t) \quad (30)$$

Where

$$s(t) = d(t)\sqrt{2S} \cos(\mathbf{w}_0 t) \quad (31)$$

This is the original BPSK signal. The property:

$$c^2(t) = 1 \quad \text{for all } t \quad (32)$$

is the key point exploited to “recover” the original BPSK signal:

$$c(t)x(t) = s(t) \quad (33)$$

If the receiver possesses a copy of the PN sequence and can synchronize the local copy with the received signal  $x(t)$ , it is able to de-spread the signal and recover the transmitted data.

### 2.2.1.1 Constant Power Broadband Noise Jammer

A jammer,  $J(t)$ , with constant power  $J$  is shown in Figure 10. The system is also assumed to have no noise from the transmission channel. An ideal BPSK demodulator is assumed after the received signal  $y(t)$  is multiplied by the PN sequence. The channel output is:

$$y(t) = x(t) + J(t) \quad (34)$$

This is multiplied by the PN sequence  $c(t)$ :

$$\begin{aligned} r(t) &= c(t)y(t) \\ &= c(t)x(t) + c(t)J(t) \\ &= s(t) + c(t)J(t) \end{aligned} \quad (35)$$

This term shows the original BPSK signal plus a noise given by  $c(t)J(t)$ . The output of the conventional BPSK detector is then:

$$r = d\sqrt{E_b} + n \quad (36)$$

where:

$d$  is the data bit for the actual  $T_b$  second interval.

$E_b = ST_b$  is the bit energy.

$n$  is the equivalent noise component.

$n$  is further defined as:

$$n = \sqrt{\frac{2}{T_b}} \int_0^{T_b} c(t)J(t) \cos(\omega_0 t) dt \quad (37)$$

The usual decision rule for BPSK is:

$$\hat{d} = \begin{cases} 1, & \text{if } r > 0 \\ -1, & \text{if } r \leq 0 \end{cases} \quad (38)$$

### 2.2.2 Coded Direct Sequence Spread Binary Phase-Shift-Keying

Several types of coding techniques can be used that provide extra gain and force the worst case jammer to be a constant power jammer. Coding techniques usually require the data rate to be decreased or the bandwidth increased because of the redundancy inherent to the coding. In spread spectrum systems, coding does not require an increase of the bandwidth or decrease of the bit rate. These properties can be seen in a simple example. If  $k=2$  (constant length) the rate is  $R=1/2$  bits per coded symbol of convolutional code. For each data bit of the sequence  $\{d\}$ , the encoder generates two coded bits. For the  $k^{th}$  transmission interval, the two data bits are:

$$a_k = (a_{k1}, a_{k2}) \quad (39)$$

where:

$$\begin{aligned} a_{k1} &= d_k \\ a_{k2} &= \begin{cases} 1 & d_k = d_{k-1} \\ -1 & d_k \neq d_{k-1} \end{cases} \end{aligned} \quad (40)$$

If  $T_b$  is the data bit time, each coded bit time is given by:

$$T_s = \frac{T_b}{2} \quad (41)$$

Defining:

$$a(t) = \begin{cases} a_{k1} & kT_b \leq t < (k+1/2)T_b \\ a_{k2} & (k+1/2)T_b \leq t < (k+1)T_b \end{cases} \quad (42)$$

$$k = \text{integer}$$

In Figure 11 the uncoded data signal  $d(t)$ , the PN sequence  $c(t)$  and the coded signal  $a(t)$  are shown for  $N=6$ . In Figure 12 the multiplied signals  $d(t)c(t)$  and  $a(t)c(t)$  are shown. With ordinary BPSK, the coded signal  $a(t)$  would have twice the bandwidth of the uncoded signal; but after spreading with the PN sequence, the final bandwidth is the same as the original. One of the simplest coding schemes is the “repeat code.” It sends  $m$  bits with the same value,  $d$ , for each data bit. The rate is then  $R=1/m$  bits per coded symbol. In this case, the resulting coded bits are:

$$a = (a_1, a_2, \dots, a_m) \quad (43)$$

Where:  $a_i = d \quad i=1,2,\dots,m$  (44)

Also, each coded bit  $a_i$  has a transmission time of:

$$T_s = \frac{T_b}{m} \quad (45)$$

It is very important to note that if  $m < N$ , the bandwidth of the spread signal does not change. The complete coded DS/BPSK system is shown in Figure 13.

The interleaver scrambles the bits in time at the transmission, and the deinterleaver reconstructs the data sequence at the receiver. After the interleaver, the signal is BPSK modulated and then multiplied by the PN sequence. At this point the transmitted DS/BPSK signal looks like the one in Eq. ( 30 ).

$$x(t) = c(t)s(t)$$

where  $s(t)$  is the common BPSK (with coding). The input at the receiver is the same as that in Eq. ( 34 ):

$$y(t) = x(t) + J(t)$$

After multiplication with  $c(t)$  (de-spreading), it becomes Eq. ( 35 ):

$$r(t) = s(t) + c(t)J(t)$$

The output of the detector after the de-interleaver is given by:

$$r_i = a_i \sqrt{\frac{E_b}{m}} + Z_i n_i \quad (46)$$

$$i = 1, 2, \dots, m$$

where  $n_1, n_2, \dots, n_m$  are independent zero mean Gaussian random variables with variance  $N_J/(2r)$ .  $r$  is the fraction of time that the pulse jammer is on, and  $Z_i$  is the jammer state:

$$Z_i = \begin{cases} 1 & \text{jammer on during } a_i \text{ transmission} \\ 0 & \text{jammer off during } a_i \text{ transmission} \end{cases} \quad (47)$$

With probability equal to:

$$\Pr\{Z_i = 1\} = r$$

$$\Pr\{Z_i = 0\} = 1 - r \quad (48)$$

### 2.2.2.1 Interleaver and Deinterleaver

The idea of using an interleaver to scramble the data bits at transmission and a deinterleaver to unscramble the bits at reception causes the pulse jamming interference on each affected data bit to be independent from each other. In the ideal interleaving and deinterleaving process, the

variables  $Z_1, Z_2, \dots, Z_m$  become independent random variables. Assume that there is no interleaver and/or deinterleaver in the system shown in Figure 13. The output of the channel is given by:

$$r_i = d \sqrt{\frac{E_b}{m}} + Z n_i \quad (49)$$

$$i = 1, 2, \dots, m$$

and because there is no interleaver/deinterleaver:

$$a_i = d$$

$$Z_i = Z \quad (50)$$

$$i = 1, 2, \dots, m$$

Also, it is assumed that the jammer was on during the whole data bit transmission  $T_b$ . Because there is no interleaver/deinterleaver, the optimum decision rule is:

$$r = \sum_{i=1}^m r_i \quad (51)$$

$$= d \sqrt{m E_b} + Z \sum_{i=1}^m n_i$$

Eq. ( 38 ) is used as a decision rule:

$$\hat{d} = \begin{cases} 1, & \text{if } r > 0 \\ -1, & \text{if } r \leq 0 \end{cases}$$

This bit error probability is the same for uncoded DS/BPSK; this means that without a interleaver/deinterleaver, there is no difference between uncoded systems and simple repeat code systems. Therefore, the use of a interleaver/deinterleaver is mandatory in order to achieve a good error probability measure against a pulse jammer.

Selection of the decision technique that determines the value of the coded bits  $\{r\}$  requires knowledge about the state of the channel. With an ideal interleaver/deinterleaver, the output of the channel is given by Eq. ( 46 ):

$$r_i = a_i \sqrt{\frac{E_b}{m}} + Z_i n_i$$

$$i = 1, 2, \dots, m$$

where  $Z_1, Z_2, \dots, Z_m$  and  $n_1, n_2, \dots, n_m$  are considered to be independent random variables. The decoder takes  $r_1, r_2, \dots, r_m$  and finds  $d_1, d_2, \dots, d_m$  with possible values of  $1$  or  $-1$ . This analysis is valid only for the instances where the state of the channel is unknown (there is no information regarding the state of the jammer signal).

### 2.2.2.2 Hard Decision Decoder

The hard decision decoder performs a binary decision on each coded bit received:

$$\hat{d}_i = \begin{cases} 1 & r_i > 0 \\ -1 & r_i \leq 0 \end{cases} \quad (52)$$

$$i = 1, 2, \dots, m$$

The final decision in decoding the transmitted bit is:

$$\hat{d}_k = \begin{cases} 1 & \sum_{i=1}^m \hat{d}_i > 0 \\ -1 & \sum_{i=1}^m \hat{d}_i \leq 0 \end{cases} \quad (53)$$

### 2.2.2.3 Interleaver Matrix

The interleaving techniques will improve the performance in pulse jammer environments because it makes the noise components become statistically independent variables. A block interleaver with depth  $I=5$  and interleaver span  $H=15$  is shown in Figure 14. The coded symbols are written to the interleaver matrix along columns, while the transmitted symbols are read out of the matrix along rows. If the coded symbol sequence is  $x_1, x_2, x_3, \dots$  the sequence that comes out of the interleaver matrix is  $x_1, x_{16}, x_{31}, x_{46}, x_{61}, \dots$ . At the receiver, the deinterleaver performs the inverse process, writing symbols into rows and reading them by columns. A jamming pulse of duration  $b$  symbols, with  $b \leq I$  will result in these jammed symbols at the deinterleaver output to be separated at least by  $H$  symbols.

## 2.3 Synchronization of Spread-Spectrum Systems

Because a pseudorandom sequence PN is used at the transmitter to modulate the signal, the first requirement at the receiver is to have a local copy of this PN sequence. The copy is needed to de-spread the incoming signal. This is done by multiplying the incoming signal by the local PN sequence copy. To accomplish a good de-spreading, the local copy has to be synchronized with the incoming signal and the PN sequence that was used in the spreading process.

The process of synchronization is usually performed in two steps: first, a coarse alignment of the PN sequence is done with a precision of less than a “chip.” This is called “PN acquisition.” After this, a fine synchronization takes care of the final alignment and corrects the small differences in the clock during transmission. This is called “PN tracking.” Theoretically, acquisition and tracking can be done in the same step with a structure of matched filters or correlators searching with high resolution the incoming signal and comparing it with the local PN sequence.

### 2.3.1 Fast Fourier Transform (FFT) Scalar Filters

These filters are implemented in the frequency domain, and they use the Fast Fourier Transform (forward and backward). They work over a set of  $N$  samples (usually in the frequency domain) [ 7 ]. The block diagram of an adaptive digital filter is shown in Figure 15.

Where:

$s(n)$  is the input signal

$n(n)$  is the noise (unwanted) signal

$r(n)$  is the input to the filter

$R(m)$  is the frequency representation of the signal ( $n$ )

$H(m)$  is the transfer function of the filter

$C(m)$  is the output (in frequency domain) after the filter is applied

$G(m)$  is the transfer function of the post-processing filter

$P(m)$  is the output after the post-processing filter

$p(n)$  is the output signal in the time domain

The following relationships are given:

$$\begin{aligned}
r(n) &= s(n) + n(n) \\
R(m) &= \text{FFT}(r(n)) \\
C(m) &= H(m)R(m) \\
P(m) &= G(m)C(m) \\
p(n) &= \text{FFT}^{-1}(P(m))
\end{aligned} \tag{54}$$

### 2.3.1.1 High-resolution Detection FFT Scalar Filter

The high-resolution detection filter outputs a peak when the desired signal  $s(n)$  and noise  $n(n)$  are applied to it. The transfer function is given by:

$$H(m) = \frac{S^*(m)}{|S(m)|^2 + |N(m)|^2} \tag{55}$$

This version of high-resolution detection assumes that the noise and the signal are uncorrelated (orthogonal). The output of this filter  $C(m)$  must be transformed to the time domain to detect the level and the position of the peak on the output vector  $c(n)$ . This position can be interpreted as the exact point where the desired signal starts within the processed set of samples  $N$ .

### 2.3.1.2 Adaptive Filtering

Adaptive filters require a learning process and use adaption techniques to form the transfer function of the desired filter  $H(m)$ . The components of the transfer function are updated periodically with actual values taken from the signal or with estimates made using stored data. The class 1/3 high-resolution detection filter is given by [ 7 ]:

$$H(m) = \frac{S^*(m)}{\langle |R(m)|^2 \rangle} \tag{56}$$

where  $S^*(m)$  is the conjugate of the spectrum of the desired signal to detect and  $|R(m)|$  is the magnitude of the spectrum of the actual input of the system.

The expression  $\langle R(m) \rangle$  is used to denote the “smoothing” process. This process is done to estimate the average spectrum of the signal plus noise from the actual input of the system. The smoothing used is called “inner block averaging” or “frequency domain averaging” and it is defined as:

$$\begin{aligned}
\langle R(j\omega) \rangle &= \frac{1}{2p} R(j\omega) * B(j\omega) \\
\text{or } r_b &= r(t)b(t)
\end{aligned} \tag{57}$$

The frequency averaging window  $B(j\omega)$  is convolved with the spectrum of the input signal. This is equivalent to a temporal weighting of the input  $r(t)$  by  $b(t)$  in the time domain. The window is usually selected to be a percentage of the input vector length.

## 3 PROPOSED SYSTEM

Different systems have been applied to watermarking of audio signals. All of them are classified as “steganographic systems” because they deal with the concept of hiding data within the signal. Boney et al. [ 23 ] proposed a system where a PN sequence was filtered using a filter that approached the masking characteristics of the human auditory system in the frequency and

time domains. Some other techniques have been imported from the fields of video and still image watermarking. Cox [ 24 ] proposes a multiplatform system capable of extract a pseudorandom sequence without the use of the original unwatermarked data.

The watermarking algorithm proposed in this paper mixes the psychoacoustic auditory model and the spread spectrum communication technique to achieve its objective. It is comprised of two main steps: first, the watermark generation and embedding and second, the watermark recovery. The watermark generation and embedding process is shown in Figure 16. A bit stream that represents the watermark information is used to generate a noise-like audio signal using a set of known parameters to control the spreading. At the same time, the audio (i.e. music) is analyzed using a psychoacoustic auditory model. The final masking threshold information is used to shape the watermark and embed it into the audio. The output is a watermarked version of the original audio that can be stored or transmitted.

The watermark recovery is shown in Figure 17. The input is the watermarked audio after transmission (i.e. music + noise, low quality, etc). An auditory psychoacoustic model is used to generate a residual. At the same time as the known parameters are used to generate the header of the watermark. Using an adaptive high-resolution filter, all the residual is scanned to find all the occurrences of the known header and therefore the initial position of each possible watermark. After this, the same known parameters used to generate the header are used to de-spread and recover the watermark.

### 3.1 WATERMARK GENERATION AND EMBEDDING

#### 3.1.1 WATERMARK GENERATION

The objective of the watermark generation is to generate a watermark audio signal  $x(t)$  that contains the watermark bit stream data. This watermark signal can be transmitted and then processed for data recovery. The technique used to generate the watermark signal  $x(t)$  is “coded DS/BPSK spread spectrum.” The process is condensed in Figure 18.

Where:

- $\{w\}$  is the original digital bit stream(watermark)
- $m$  is the repetition code factor
- $\{w_R\}$  is the watermark after the coding process (repeat code)
- $I,H$  = width and length of the interleaver matrix
- $\{w_I\}$  is the watermark after the interleaver process
- $\{header\}$  = is the header sequence
- $\{d\} = \{header\} + \{w_I\}$  = sequence to be spread and transmitted
- $f_0$  = frequency used by the BPSK modulator

The process can be explained with a simple example: Let  $\{w\}$  be the watermark bit stream. All the bit streams used are bipolar (value 1 or -1). Defining  $\{w\}$  with a length of 16 bits as the sequence:

$$\{w\} = \{ 1 \quad 1 \quad -1 \quad 1 \quad -1 \quad -1 \quad 1 \quad -1 \mid 1 \quad 1 \quad -1 \quad 1 \quad 1 \quad 1 \quad -1 \quad -1 \}$$

Using Eq. ( 43 ) to generate the repeat code, and choosing  $m=3$ , the  $\{w_R\}$  sequence is:

$$\{w_R\} = \left\{ \begin{array}{cccccccccccc} 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & 1 & 1 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & -1 & -1 & -1 \\ 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & -1 \end{array} \right\}$$

The next step is to perform interleaving. To do this, the values of the interleaving matrix are chosen; in this case,  $I=5$ ,  $H=10$ , (see Figure 14). The resulting matrix is shown in Figure 19. The last two spaces are padded with 1's. Using the interleaving matrix, the output sequence  $\{w_I\}$  is:

$$\{w_I\} = \left\{ \begin{array}{cccccccccccc} 1 & 1 & 1 & -1 & 1 & 1 & 1 & -1 & -1 & 1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 & -1 & 1 & 1 \\ -1 & 1 & -1 & 1 & 1 & -1 & -1 & -1 & 1 & 1 & -1 & -1 \\ -1 & 1 & 1 & -1 & -1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & & & & & & & & & & \end{array} \right\}$$

The selected header is a sequence usually composed by 1's.

$$\{header\} = \{1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1\}$$

The final data sequence  $\{d\}$  is obtained concatenating the  $\{header\}$  and the  $\{w_I\}$ :

$$\{d\} = \{header\} + \{w_I\}$$

$$\{d\} = \left\{ \begin{array}{cccccccccccc} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & 1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 & -1 \\ -1 & 1 & -1 & -1 & 1 & -1 & 1 & -1 & 1 & 1 & -1 & 1 \\ -1 & 1 & 1 & -1 & -1 & -1 & 1 & 1 & -1 & -1 & -1 & 1 \\ 1 & -1 & -1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{array} \right\}$$

The PN sequence  $\{c\}$  can be generated by any means. Usually this is done using a pseudorandom number generator. In this case, the PN sequence is assumed to be long enough to spread a complete bit stream (header and data) without repeating any portion of it. The important factor is that the transmitter and the receiver must have a copy of the whole PN sequence  $\{c\}$ .

This sequence is ideally uncorrelated with the  $\{d\}$  sequence, and has the form:

$$\{c\} = \{1 \ -1 \ 1 \ 1 \ -1 \ -1 \ 1 \ -1 \ 1 \ 1 \ -1 \ 1 \ \dots\}$$

### 3.1.1.1 Spread Spectrum Parameter Selection

Audio signals are usually considered to be baseband signals [ 21 ]. The described spread spectrum technique can be applied to passband systems (with  $f_0 > 0$ ) or baseband systems ( $f_0 = 0$ ) without losing generality. The selection of all the parameters is based on the considerations of how the overall watermarked audio signal will be transmitted or stored. The frequency response of those systems determines which frequencies are likely to be present at the receiver. Let a baseband bandlimited signal, with no modulation ( $f_0 = 0$ ) have the magnitude spectrum shown in Figure 20. With amplitude modulation ( $f_0 > 0$ ), the spectrum will have the form shown in Figure 21. FS is the sampling frequency of the system. To avoid aliasing because of the use of modulation, the modulation frequency should be:

$$Rc \leq f_0 \leq \frac{FS}{2} - Rc \quad (58)$$

If a system possesses a lower frequency limit  $LF$  and/or an upper frequency limit  $HF$ , the modulation frequency  $f_0$  have to be selected in a way that the sidebands fall between the lower an upper limits, as shown in Figure 22. If a sideband falls outside of these limits, aliasing or data loss could result. Taking into account, the selection of parameters should be done using:

$$LF + Rc \leq f_0 \leq HF - Rc$$

$$LF \geq 0, HF \leq \frac{FS}{2} \quad (59)$$

The parameters selected must satisfy Eq. ( 58 ) and Eq. ( 59 ), along with the following relationships:

- $R_d$  = is the data bits per second
- $m$  = is the repetition code factor
- $N$  = is the spreading factor, Eq. ( 26 )
- $R_b = R_d * m$  is the coded bits per second
- $T_b = 1/R_b$  is the time of each coded bit
- $R_c = N * R_b$  is the PN sequence bits per second
- $T_c = T_b / N$  is the time of each PN bit or “chip”

Assuming a frequency response similar to FM Radio [ 22 ] with  $LF = 50 \text{ Hz}$  and  $HF = 15000 \text{ Hz}$ , for the actual example, a set of spread spectrum parameters that satisfy all the requirements is:

- $N = 3$
- $m = 3$
- $R_d = 100 \text{ bits/sec}$
- $R_b = 300 \text{ bits/sec}$
- $R_c = 900 \text{ bits/sec}$
- $f_0 = 3500 \text{ Hz}$

Note that  $N$  and  $m$  are selected with small values for this example. The modulation is done using Eq. ( 31 ):

$$s(t) = d(t) \sqrt{2S} \cos(\omega_0 t)$$

The spreading is done using Eq. ( 30 ):

$$x(t) = c(t)s(t)$$

The output of the system is the watermarked audio waveform  $x(t)$  shown in Figure 23.

### 3.1.2 FRAME SEGMENTATION

To overcome the potential problem of the audio signal or the watermark signal being too long to be processed using a single FFT, the signal is segmented in short overlapping segments, processed and added back together [ 8 ]. Another consideration for the watermark algorithm is that the audio signal has to be longer than the watermark signal. Therefore, the watermark can be repeated several times during the duration of the audio signal. This redundancy is one of the important features in the watermarking algorithm. Figure 24 shows audio and watermark signals that will be segmented. The watermark is repeated several times. If the total length of the audio signal is  $LENGTH$  samples, the desired length of the analysis frame is  $BLOCK$  samples, and the

overlap between consecutive frames is  $OVERLAP$  samples, the total number of  $FRAMES$  is given by:

$$FRAMES = \frac{LENGTH - OVERLAP}{BLOCK - OVERLAP} \quad (60)$$

In Figure 24 two equal length frames were selected to be processed. One from the audio signal and the other from the respective point in the watermark signal. The last frame is zero-padded if it is shorter than  $BLOCK$  samples. These padded samples are discarded in post processing. From this point on, all processes described are applied to the audio or watermark signal frames, not the entire signal.

### 3.1.3 FREQUENCY REPRESENTATION

The Short Time Fourier Transform (STFT) is used to acquire a frequency representation of the actual frames. Before doing the STFT, a Hamming window is applied to both signals [ 7 ], [ 8 ]. This improves the representation of the signal in the frequency domain reducing the leakage. If  $s(t)$  is the actual audio signal frame and  $x(t)$  the actual watermark signal frame, then the windowing is done using:

$$sw(t) = s(t)w(t) \quad (61)$$

$$xw(t) = x(t)w(t) \quad (62)$$

The Hamming window is defined as:

$$w(n) = 0.54 + 0.46 \cos\left(\frac{2n\pi}{BLOCK}\right)$$

$$n = 1, 2, \dots, BLOCK \quad (63)$$

$$w(t) = w(nT)$$

$$T = \text{sampling period}$$

The frequency representation of the audio frame is:

$$Sw(j\omega) = FT\{sw(t)\} \quad (64)$$

and the watermark frame:

$$Xw(j\omega) = FT\{xw(t)\} \quad (65)$$

The power spectra is found using Eq. ( 3 ):

$$Sp(j\omega) = |Sw(j\omega)|^2 \quad (66)$$

The indices of the actual frequency representations have to be mapped to the Bark scale. Once this index mapping is done, the representation in the critical band scale is formed by mapping the components to the respective position on the critical band axis. The relationship between each component index,  $i$ , and the corresponding frequency,  $f_i$ , that it represents is given by:

$$f_i = \frac{(i-1) * FS}{BLOCK}$$

$$i = 1, 2, \dots, \frac{BLOCK}{2} \quad (67)$$

$$FS = \text{Sampling Frequency}$$

The relationship between each frequency  $f_i$  and the bark scale or critical band scale  $z_I$  is found using Eq. ( 1 ):

$$z_i = 13 \tan^{-1} \left( \frac{0.76 * f_i}{1000} \right) + 3.5 \tan^{-1} \left( \left( \frac{f_i}{7500} \right)^2 \right)$$

This relationship between each component index  $i$  and the frequency  $f_i$  or critical band  $z_i$  that it represents can be calculated at the beginning of the algorithm and stored in a table. The energy per critical band is calculated using Eq. ( 4 ):

$$Spz(z) = \sum_{w=LBZ}^{HBZ} Sp(jw)$$

Where:  $z = 1, 2, \dots$ , total of critical bands  $Zt$ ;  $LBZ$  and  $HBZ$  the lower and higher frequencies in the critical band  $z$ . Figure 25 (a) shows the original audio frame  $s(t)$  in the time domain and the shape of the Hamming window  $w(t)$ ; (b) shows the  $sw(t)$  frame after the windowing process; (c) shows the magnitude of  $Sw(jw)$ , and (d) shows the power spectrum  $Sp(jw)$  and the energy per critical band  $Spz(z)$ .

### 3.1.4 BASILAR MEMBRANE SPREADING FUNCTION

The basilar membrane spreading function determines how much of the energy of each critical band is contributed to the neighboring bands. The spreading function  $B(z)$  is calculated using Eq. ( 5 ):

$$B_k = 15.91 + 7.5(k + 0.474) - 17.5\sqrt{1 + (k + 0.474)^2}$$

$$k = \dots - 2, -1, 0, 1, 2, \dots$$

The spreading across bands is computed by the convolution of the spreading function  $B(z)$  and the energy per critical band  $Spz(z)$ , using Eq. ( 6 ):

$$Sm(z) = Spz(z) * B(z)$$

Figure 26 (a) shows the energy per critical band  $Spz(z)$ , (b) shows the spreading function  $B(z)$  for 9 points, and (c) shows the spread energy per critical band  $Sm(z)$ .

### 3.1.5 MASKING THRESHOLD ESTIMATE

The Spectral Flatness Measure (SFM) of the actual audio frame  $Sw(jw)$  is taken using Eq. ( 7 ):

$$SFM_{dB} = 10 \log_{10} \left\{ \frac{\prod_{z=1}^{Zt} Spz(z)}{\frac{1}{Zt} \sum_{z=1}^{Zt} Spz(z)} \right\}^{\frac{1}{Zt}}$$

with  $Zt =$  total number of critical bands in each frame

The energy per critical band  $Spz(z)$  is used rather than spread energy per critical band  $Sm(z)$  to avoid false results due to smoothing of the signal. The tonality factor  $\mathbf{a}$  is then calculated using Eq. ( 8 ):

$$\mathbf{a} = \min \left( \frac{SFM_{dB}}{SFM_{dB \max}}, 1 \right)$$

with  $SFM_{dB\max} = -60dB$ .

The masking energy offset  $O(z)$  is then calculated using Eq. ( 9 ):

$$O(z) = \mathbf{a}(14.5 + z) + (1 - \mathbf{a})5.5$$

The raw masking threshold,  $Traw(z)$ , is calculated with Eq. ( 10 ):

$$Traw(z) = 10^{\left(\log_{10}(Sm(z)) - \frac{O(z)}{10}\right)}$$

The raw masking threshold is normalized using Eq. ( 11 ):

$$Tnorm(z) = \frac{Traw(z)}{P_z}$$

where:  $P_z$  = number of points in each band  $z$   
 $z = 1, 2, \dots, Zt$

To calculate the final masking threshold  $T$  it is necessary to first calculate the hearing threshold (or threshold in quiet)  $TH$ . It is defined as a sinusoidal tone of 4000 Hz with one bit of dynamic range. Using Eq. ( 12 ):

$$TH = \max(|Pp(j\mathbf{w})|)$$

Where:  $Pp(j\mathbf{w})$  = power spectrum of the probe signal  $p(t)$   
 $p(t) = \sin(2\mathbf{p}4000t)$

Then the final masking threshold  $T$  is calculated using Eq.( 13 ):

$$T(z) = \max(Tnorm(z), TH)$$

with:  $z=1, 2, \dots, Zt$

Figure 27 (a) shows the raw masking threshold  $Traw(z)$  and (b) shows the normalized threshold  $Tnorm(z)$ .

### 3.1.6 WATERMARK SPECTRAL SHAPING

The final masking threshold  $T$  is used to determine which components of the audio signal  $Sw(j\mathbf{w})$  can be removed without affecting the perceptual quality of the signal. The power spectrum  $Sp(j\mathbf{w})$  is compared against the final masking threshold  $T$ . The components that fall below it are removed in  $Sw(j\mathbf{w})$ . The new frame with only the components above the threshold is called  $Swnew(j\mathbf{w})$ . Eq. ( 14 ) is used:

$$Swnew_i(j\mathbf{w}) = \begin{cases} Sw_i(j\mathbf{w}) & Sp_i(j\mathbf{w}) \geq T(z) \\ 0 & Sp_i(j\mathbf{w}) < T(z) \end{cases}$$

$i = 1, 2, \dots$  number of components  
 $z, \mathbf{w}$  according to component  $i$

Then the unneeded components of the watermark signal  $Xw(j\mathbf{w})$  are removed. These components correspond to the non-removed components in  $Sw(j\mathbf{w})$ . Eq. ( 15 ) is used:

$$Xwnew_i(j\mathbf{w}) = \begin{cases} 0 & Sp_i(j\mathbf{w}) \geq T(z) \\ Xw_i(j\mathbf{w}) & Sp_i(j\mathbf{w}) < T(z) \end{cases}$$

$i = 1, 2, \dots$  number of components  
 $z, \mathbf{w}$  according to component  $i$

The factors that will shape the new watermark  $X_{wnew}(j\omega)$  are found using Eq. ( 18 ):

$$F_z = A \frac{\sqrt{T(j\omega)}}{\max(|X_{wnew}(j\omega)|)}$$

$$z = 1, 2 \dots Z_t$$

$$\omega = LBZ \text{ to } HBZ \text{ for each band } z$$

The square root of the final threshold is divided by the maximum magnitude component found in the energy of the new watermark in each critical band. Each one of these factors is scaled using the gain  $A$ , that varies from 0 to 1, and controls the overall magnitude of the watermark signal in relation with the audio signal.

Each one of the components in each critical band  $k$  is scaled by the corresponding factor using Eq. ( 19 ):

$$X_{final}(j\omega) = X_{wnew}(j\omega)F_z$$

$$z = 1, 2 \dots Z_t$$

$$\omega = LBZ \text{ to } HBZ \text{ for each band } z$$

Figure 28 shows the final masking threshold and the watermark signal before shaping (a) and after shaping (b). Note that the watermark falls below the threshold of masking. The factor  $A$  gives control of “how much gain” will have the watermark related with the masking threshold ( $A$  is a value from 0 to 1).

### 3.1.7 AUDIO AND WATERMARK SIGNAL COMBINATION

The final output  $OUT(j\omega)$  is the sum of the new audio,  $S_{wnew}(j\omega)$ , and the final watermark  $X_{final}(j\omega)$ . This is given by the Eq. ( 20 ):

$$OUT(j\omega) = S_{wnew}(j\omega) + X_{final}(j\omega)$$

Figure 29 shows the final masking threshold  $T_{final}(z)$ , and the power spectrum of (a)  $S_{wnew}(j\omega)$ , (b)  $X_{final}(j\omega)$ , and (c)  $OUT(j\omega)$ .

### 3.1.8 TRANSFORMATION TO THE TIME DOMAIN

The Inverse Fourier Transform is used to convert the frequency domain information back to the time domain.

$$out(t) = \text{IFT}\{OUT(j\omega)\}$$

This output frame  $out(t)$  is added to the correspondent point at the total time domain output  $output(t)$ . The next frames of audio and watermark signals are taken, and the process is repeated.

## 3.2 DATA RECOVERY

The watermarked audio signal is intended to be transmitted through a diverse number of channels. In some cases, the channel will introduce noise, convert several times from digital to analog and analog to digital, or even use a psychoacoustic auditory model to process the audio signal. The watermark bit stream should survive the transmission and be recoverable.

A very important characteristic is that the developed system does not require access to the original audio signal (before watermark) to extract the watermark at the receiving. The process of recovery uses the psychoacoustic auditory model, but in this case the goal is to remove all the audio components that have less probability of belonging to the watermark signal. This means

that the masking threshold is calculated and the components above it are removed. The final signal is the “residual.” This residual is then analyzed to find the possible points where the watermark is present. If some criterion is applied, the majority of the false points detected can be eliminated (i.e. rejecting points too close to fit a watermark). Synchronization and recovery of the watermark bit stream are then performed.

### 3.2.1 MASKING THRESHOLD AND RESIDUAL SIGNAL

The watermarked audio signal after the transmission is symbolized as  $s_2(t)$ . The process described in sections 3.1.2 to 3.1.5 is used to calculate the frames  $sw_2(t)$ , frequency representation  $Sw_2(j\omega)$ , and masking threshold  $T_2$ , respectively. The residual signal  $R(j\omega)$  is defined as the signal composed of the components below the masking threshold. Eq. ( 14 ) can be changed to:

$$R_i(j\omega) = \begin{cases} Sw_{2i}(j\omega) & Sp_{2i}(j\omega) \leq T_2(z) \\ 0 & Sp_{2i}(j\omega) > T_2(z) \end{cases} \quad (68)$$

$i = 1, 2, \dots$  number of components  
 $z, \omega$  according to component  $i$

### 3.2.2 RESIDUAL EQUALIZATION

The spectrum of the residual  $R(j\omega)$  is then shaped to be flat. Eq. ( 18 ) can be modified to shape all the maximum components of each band to be at equal levels. The factors are found using:

$$F_z = \frac{1}{\max(|R(j\omega)|)} \quad (69)$$

$z = 1, 2, \dots Zt$   
 $\omega = LBZ$  to  $HBZ$  for each band  $z$

Each one of the components in each critical band  $z$  is scaled by the corresponding factor  $F_z$  using Eq. ( 19 ):

$$R_{final}(j\omega) = R(j\omega)F_z$$

$z = 1, 2, \dots Zt$   
 $\omega = LBZ$  to  $HBZ$  for each band  $z$

### 3.2.3 TIME DOMAIN RESIDUAL

The residual is taken back to the time domain using the Inverse Fourier Transform IFT.

$$r(t) = \text{IFT}\{R_{final}(j\omega)\}$$

The time domain  $r(t)$  frame is added to the total time domain residual signal  $residual(t)$  at the point specified by the frame segmentation step. The next frame is then processed.

### 3.2.4 SYNCHRONIZATION WITH WATERMARK HEADER

To be able to synchronize and to have a good de-spreading of the watermark signal, it is necessary to have knowledge of the parameters used at the generation of the watermark signal, such as  $f_0$ ,  $T_b$ ,  $m$ ,  $H$ ,  $I$ ,  $N$ ,  $\{header\}$ ,  $\{c\}$ , etc.

### 3.2.4.1 *header(t)* Signal Generation

The first step is to generate a *header(t)* waveform signal using the process of section 3.1.1, except that only the  $\{header\}$  sequence is used as the input sequence. This audio signal will be used to locate the exact positions of the watermark signals in the *residual(t)* signal. Frame segmentation as explained in section 3.1.2 is also required in order to analyze the whole *residual(t)* signal. The parameters for the frame segmentation are chosen to have up to two *header(t)* signals in each frame. Therefore, BLOCK is equal to twice the number of samples in *header(t)*, and OVERLAP is equal to one half the number of samples in *header(t)*. The resulting frame taken from *residual(t)* with BLOCK length is called *r(t)*.

### 3.2.4.2 *header(t)* Position Detection

Eq. ( 56 ) describes an adaptive high-resolution filter that can be used to detect the presence of *header(t)* in the *r(t)* frame and therefore, all the occurrences of *header(t)* in the *residual(t)* audio signal.

$$H(j\omega) = \frac{HEADER^*(j\omega)}{\langle |R(j\omega)|^2 \rangle}$$

Where:

$$R(j\omega) = \text{FFT}(r(t))$$

$$HEADER(j\omega) = \text{FFT}(header(t))$$

The denominator of the filter is the smoothed version of  $|R(j\omega)|^2$ . Smoothing is done using Eq. ( 57 ), where  $w(t)$  is a Hanning window of width 10%. The output of the filter applied to  $R(j\omega)$  is:

$$DET(j\omega) = R(j\omega) \frac{HEADER^*(j\omega)}{\langle |R(j\omega)|^2 \rangle}$$

This result is transformed to the time domain to be analyzed.

$$det(t) = \text{real}(\text{IFFT}(DET(j\omega)))$$

A typical output of the filter, *det(t)*, is shown in Figure 30. The peak shows the position in samples where the *header(t)* signal starts in the frame *r(t)*. This detection is done for all the frames in the *residual(t)* signal, and all the positions of the peaks are stored for further analysis. A proposed criterion of analysis is to determine the minimum distance between peaks to decide which ones have more probability to represent the start of a watermark signal.

### 3.2.5 WATERMARK DE-SPREADING

For each peak position found in the *residual(t)*, a selected frame *y(t)* with the same length as the watermark signal is processed. This process is shown in Figure 31. Using Eq. ( 35 ):

$$r(t) = c(t)y(t)$$

Demodulation is performed using Eq. ( 31 ):

$$g(t) = r(t) \sqrt{\frac{2}{T_b}} \cos(2\mathbf{p}f_0 t)$$

To estimate the bit stream:

$$r_i = \int_{(i-1)T_s}^{iT_s} g(t)dt \quad (70)$$

$i = 1, 2, \dots$  total bits in bit stream

The decision rule, to form a recovered bit stream  $\{\hat{d}\}$ , is given by Eq. ( 38 ),

$$\hat{d}_i = \begin{cases} 1, & \text{if } r_i > 0 \\ -1, & \text{if } r_i \leq 0 \end{cases}$$

$i = 1, 2, \dots$  total bits in bit stream

After this decision, the  $\{header\}$  sequence is discarded from the  $\{\hat{d}\}$  bit stream. This produces the bit stream,  $\{\hat{w}_i\}$ .

### 3.2.6 WATERMARK DE-INTERLEAVING AND DECODING

The de-interleaving process is done using the same matrix used in the watermark generation in section 3.1.1 and shown in Figure 14. The bits are written into rows and read by columns to accomplish the de-interleaving process. The de-interleaved sequence is called  $\{\hat{w}_R\}$ . The decoding of the repeat code of value  $m$  is done using Eq. ( 53 ):

$$\hat{w}_k = \begin{cases} 1 & \sum_{i=1}^m \hat{w}_{Ri} > 0 \\ -1 & \sum_{i=1}^m \hat{w}_{Ri} \leq 0 \end{cases}$$

$k = 1, 2, \dots$  total bits in data sequence

The final recovered sequence  $\{\hat{w}\}$  is the recovered watermark.

## 4 SYSTEM PERFORMANCE

### 4.1 SURVIVAL OVER DIFFERENT CHANNELS

A watermarking system was implemented using a well known mathematical software package. The system was composed of two modules: watermark generation and embedding, and watermark recovery. The watermark was first generated and embedded in an audio signal. The watermarked signal was then tested for recovery of the watermark after transmission by different channels, such as sub-band encoding, digital to analog – analog to digital conversions and radio transmission.

The music used was a 26 second excerpt of the song “In the Midnight Hour” (W. Pickett & S. Cropper) performed by *The Commitments*. A sampling frequency of 44.1 KHz was used. Each of the watermarked audio signals was labeled to reflect the level of the watermark below the masking threshold (the A value), i.e. W2, W4, W6 and W8. With these parameters, a total of 35 watermarks were embedded during the duration of each signal. The four watermarked music signals and the original signal were recorded digitally on a compact disc. The computer was also equipped with a full duplex sound card with D/A A/D converters. All the radio systems were simulated using a multiplex stereo modulator, FM/AM signal generator, and ordinary consumer CD player and FM/AM radio receiver. The percentage of correct bits recovered per watermark was measured before and after transmission. Two examples are shown in Figure 32 and Figure 33. The percentage of correct bits before transmission is the continuous line, and the percentage

of correct bit after transmission is the dotted line. Also, the offset from the expected starting point of each watermark after transmission is measured (in samples), as well as the total of watermarks recovered and the average recovery percentage.

## 4.2 LISTENING TEST

One of the requirements of the watermarking system is to retain the perceptual quality of the signal. This is often referred to as “transparency.” The transparency of the watermarking algorithm was tested using three of the four watermarked audio signals (W2, W4 and W6) used in section 4.1. An ABX listening test was used as the testing mechanism. In an ABX test the listener can hear selection A (in this case the non-watermarked audio), selection B (the watermarked audio) and X (either the watermarked or non-watermarked audio). The listener is then asked to decide if selection X is equal to A or B. The number of correct answers is the basis to decide if the watermarked audio is perceptually different than the original audio and would, therefore, declare the watermarking algorithm as “non-transparent.” In the other case, if the watermarked audio is perceptually equal to the original audio, the watermarking algorithm will be declared as “transparent.”

Using the theory explained in Burstein [ 19 ], [ 20 ], different parameters were selected to find an appropriate sample size. A criterion of significance  $\alpha'=0.1$  is selected (also known as Error Type 1). The Type 2 error risk is assumed  $\beta'=0.1$ . The probability  $p1$  that a listener finds the right answer by chance is 0.5 in an ABX system. The effect size is selected as  $p2=0.7$ . With these parameters, the approximated required sample size that meets the specifications is 37.61 samples. The sample size is selected as  $n=40$ . (40 listeners per ABX set). The critical  $c$  ( $c'$ ) is the minimum number of correct samples which, together with  $n$  and  $p1$ , can produce a significance level  $\alpha$  equal to or less than the specified criterion of significance  $\alpha'$ . The calculated  $c'$  is 24.55 and can be rounded off to 25. This is the minimum number of correct answers to accept the hypothesis that the listener perceives differences between audio A and B. With  $c'=25$ , the criterion of significance becomes  $\alpha'=0.78$ , which is below the required level. The type 2 error risk  $\beta'=1.11$  and does not exceed desired level. The results and their approximate significance level are shown in Table 1.

	Sample Size	Correct Identifications	$\alpha$
W2	40	24	0.14
W4	40	19	0.50
W6	40	19	0.50

Table 1. Listening test results

## 4.3 DISCUSSION

The survival over different channels showed that after encoding, not all the watermarks could be recovered with 100% accuracy. This occurs because of the multiple factors that affect the quality of the embedded watermark, such as: the number of audio components replaced, the gain of the watermark, and the masking threshold. It is important to note that in some frames the watermark information can be very weak, even null. The spread spectrum technique employed can partially solve these problems, but if many consecutive frames have no watermark information, that specific watermark can not be recovered.

The theoretical position of the watermark and the offset of the actual watermark represent the starting position of the *{header}* of each watermark. This position will not affect the recovery of

the watermark because each watermark is embedded independently of the others. In the actual tests three different cases are seen: almost no offset, linearly increasing offset and varying offset. When no offset is seen, the original signal and the recorded signal after transmission were played at the same speed. In the cases where the offset is linearly, it is assumed that the speed of the playback device (in this case an ordinary consumer CD player was different (slightly slower) than the recording device. The last case shows the unstable speed variations of the tape device. If the speed of the playback device is close enough to the original speed, the de-spreading can be successful because the difference in alignment between the watermarked audio and the de-spreading signals (PN sequence, demodulator and *{header}*) will not greatly affect the final result.

Finally, the percentage of correct bits recovered measures quality of the recovery for each watermark. Notice that not all the watermarks are recovered (*%bits* = 0.0), and not all the watermarks are recovered in their totality but many of them were recovered with more than 80% of the bits. A good bit error detection/correction algorithm or averaging technique could substantially improve the recovery of the watermark. A very strong point in the watermarking system is the redundancy of watermarks embedded into the audio stream. In this case, each watermark lasts approximately 600 ms. Even if just a few watermarks are recovered, the goal of transmitting the watermark information within the audio signal and recovering it afterwards is accomplished.

The listening test showed that the watermark at  $-2\text{dB}$  below the masking threshold ( $W_2$ ) is the most likely to be heard, but it can not be ensured that people actually noticed the difference. For all the other watermarked signals, the results show that the process is “transparent.”

## 5 CONCLUSIONS

The proposed digital watermarking method for audio signals is based on a psychoacoustic auditory model to shape an audio watermark signal that is generated using spread spectrum techniques. The method retains the perceptual quality of the audio signal, while being resistant to diverse removal attacks, either intentional or unintentional. The recovery of the watermark is accomplished without knowledge of the original audio signal. The only information used includes the watermarked audio signal, and the parameters used for the watermark generation.

The psychoacoustic auditory model retrieves the necessary information about the masking threshold of the input audio signal. This model is a good approach that can be used for several applications such: perceptual coding, masking analysis, or watermark embedding. The spread spectrum theory describes two important Direct Sequence techniques, but the employed technique is Coded Direct-Sequence Spread Binary Phase-Shift-Keying (coded DS/BPSK). Because the normal literature about this topic is reserved for communication theory, some assumptions were made to use the theory in an audio bandwidth environment. Specifically in this case, the audio information was considered the “noise” or “jammer” signal that interferes with the watermark.

Future research could be performed in different aspects of this proposed algorithm such as:

- System performance with different types of music.
- Experimenting with different spread spectrum encoding parameters.
- Changes in the playback speed of the signal.
- Crosstalk interference.
- Multiple watermark embedding.

- Use of techniques to enhance recovery of the watermark (i.e., bit error detection/ correction, averaging, etc).
- Real - time implementation.
- Investigate different signal schemes for the generation of the PN sequence.

## 6 ACKNOWLEDGMENT

The author wishes to thank Professors Ken Pohlmann and Will Pirkle from the Music Engineering program at University of Miami for their valuable advises and feedback. Also to the Music Engineer Alex Souppa for his help as technical editor and english corrector of the author's master thesis.

## 7 REFERENCES

- [ 1 ] E. Zwicker and U. T. Zwicker, "Audio Engineering and Psychoacoustics: Matching Signals to the Final Receiver, the Human Auditory System," *J. Audio Eng. Soc.*, vol. 39, pp. 115 -126 (1991 March)
- [ 2 ] T. Sporer and K. Brandenburg, "Constraints of Filter Banks Used for Perceptual Measurement," *J. Audio Eng. Soc.*, vol. 43, pp. 107 - 115 (1995 March)
- [ 3 ] J. Mourjopoulos and D. Tsoukalas, "Neural Network Mapping to Subjective Spectra of Music Sounds," *J. Audio Eng. Soc.*, vol. 40, pp. 253 - 259 (1992 April)
- [ 4 ] J. D. Johnston, "Transform Coding of Audio Signals Using Perceptual Noise Criteria," *IEEE Journal on Selected Areas in Communications*, vol. 6, pp. 314 – 323 (1988 Feb.)
- [ 5 ] M. K. Simon, J. K Omura, R A. Scholtz and B. K. Levitt, *Spread Spectrum Communications Handbook* (McGraw-Hill, New York, 1994)
- [ 6 ] R. L. Pickholtz, D. L. Schilling, and L. B. Milstein, "Theory of Spread-Spectrum Communications – A Tutorial," *IEEE Transactions on Communications*, vol. COM-30, pp. 855 – 884 (1982 May)
- [ 7 ] C. S. Lindquist, *Adaptive & Digital Signal Processing with Digital Filtering Applications* (Steward & Sons, Miami, 1989)
- [ 8 ] L. R. Rabiner, and R. W. Schafer, *Digital Processing of Speech Signals* (Prentice Hall, New Jersey, 1978)
- [ 9 ] E. Zwicker, and h. Fastl, *Psychoacoustics Facts and Models* (Springer-Verlag, Berlin, 1990)
- [ 10 ] D. L. Nicholson, *Spread Spectrum Signal Design. LPE & AJ Systems* (Computer Science Press, Rockville, Maryland, 1988)
- [ 11 ] J.C. Neubauer and J. Herre, "Digital Watermarking and Its Influence on Audio Quality," presented at the 105<sup>th</sup> Convention of the Audio Engineering Society, *J. Audio Eng. Soc. (Abstracts)*, vol. 46, pp. 1041 (1998 November), preprint 4823.
- [ 12 ] J. G. Roederer, *The Physics and Psychophysics of Music* (Springer-Verlag, New York, 1995)
- [ 13 ] J. G. Beerends and J. A. Stemerdink, "A Perceptual Speech-Quality Measure Based on a Psychoacoustic Sound Representation," *J. Audio Eng. Soc.*, vol. 42, pp. 115 - 123 (1994 March)
- [ 14 ] J. G. Beerends and J. A. Stemerdink, "A Perceptual Audio Quality Measure Based on a Psychoacoustic Sound Representation," *J. Audio Eng. Soc.*, vol. 40, pp. 963 - 978 (1992 December)
- [ 15 ] C. Colomes, M. Lever, J. B. Rault, Y. F. Dehery and G. Faucon, "A Perceptual Model Applied to Audio Bit-Rate Reduction," *J. Audio Eng. Soc.*, vol. 43, pp. 233 - 239 (1995 April)

- [ 16 ] T. Sporer, G. Gbur, J. Herre and R. Kapust, "Evaluating a Measurement System," *J. Audio Eng. Soc.*, vol. 43, pp. 353 - 362 (1995 May)
- [ 17 ] M. R. Schroeder, B. S. Atal and J. L. Hall, "Optimizing Digital Speech Coders by Exploiting Masking Properties of the Human Ear," *J. Acoust. Soc. Am.*, vol. 66, pp. 1647 – 1652 (1979 Dec.)
- [ 18 ] B. Paillard, P. Mabillean, S. Morissette and J. Soumagne, "PERCEVAL: Perceptual Evaluation of the Quality of Audio Signals," *J. Audio Eng. Soc.*, vol. 40, pp. 21 - 31 (1992 Jan./Feb.)
- [ 19 ] H. Burstein, "By the Numbers," *Audio*, vol. 74, pp. 43 – 48 (1990 Feb.)
- [ 20 ] H. Burstein, "Approximation Formulas for Error Risk and Sample Size in ABX Testing," *J. Audio Eng. Soc.*, vol. 36, pp. 879 - 883 (1988 Nov.)
- [ 21 ] S. Haykin, *Communication Systems 3<sup>rd</sup> ed.* (Wiley, New York, 1994)
- [ 22 ] R. L. Shrader, *Electronic Communication 5<sup>th</sup> ed.* (McGraw Hill, New York, 1985)
- [ 23 ] L. Boney, A. H. Tewfik and K. N. Hamdy, "Digital Watermarks for Audio Signals," *IEEE Int. Conf. on Multimedia Computing and Systems*, Hiroshima, Japan (June 1996)
- [ 24 ] I. J. Cox, "Spread Spectrum Watermark for Embedded Signalling", *United States Patent 5,848,155* (1998 Dec)

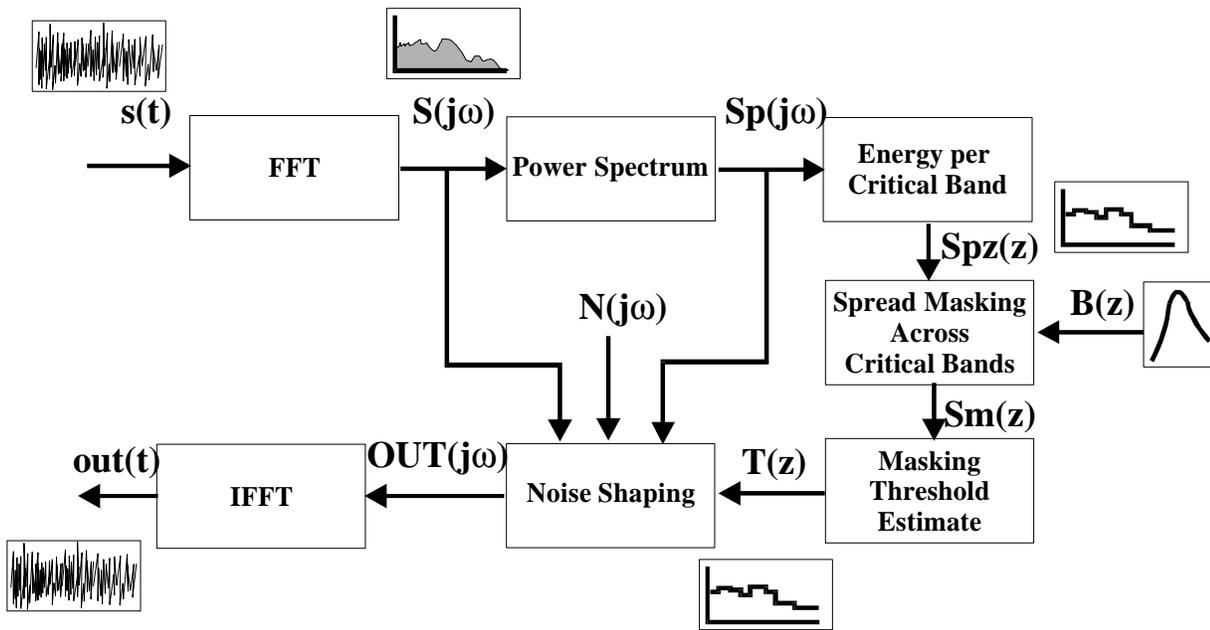


Figure 1. Psychoacoustic auditory model

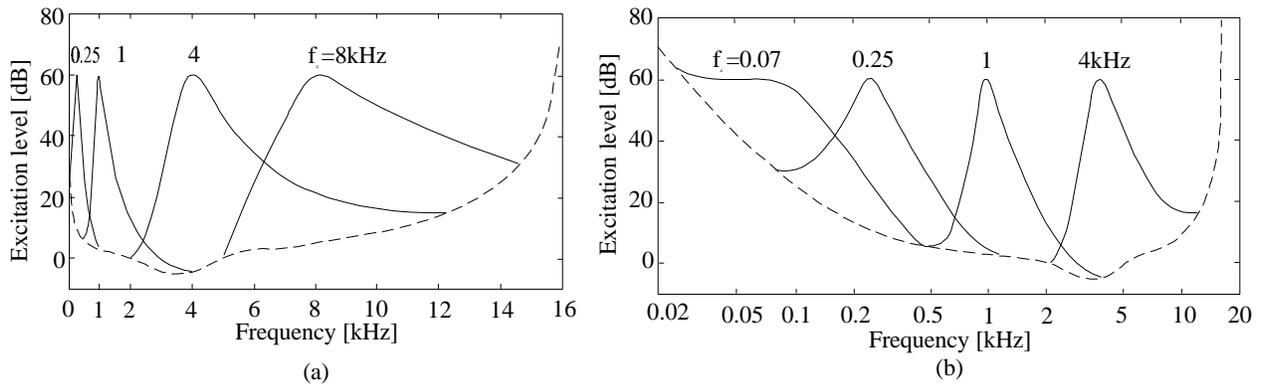


Figure 2. Masking curves in (a) linear and (b) logarithmic frequency scale [ 1 ]

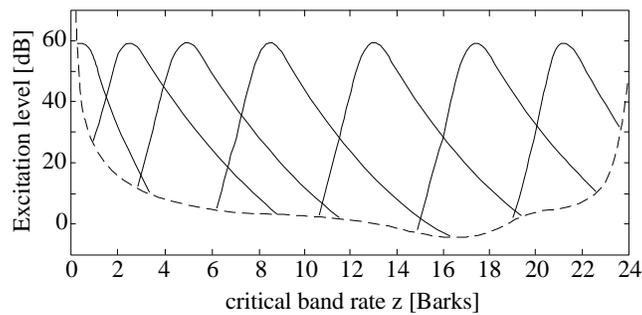


Figure 3. Excitation level versus critical band rate for narrow band noises with various center frequencies [ 1 ]

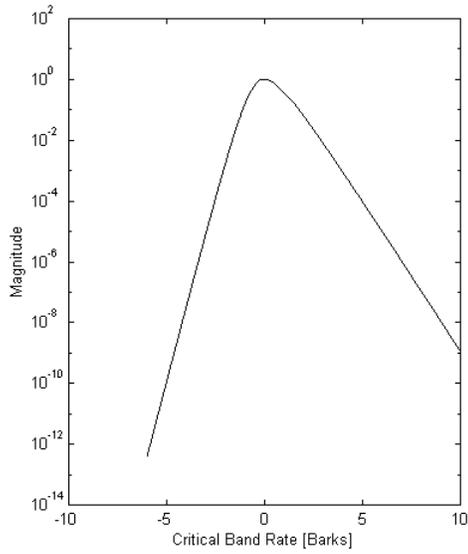


Figure 4. Model of the spreading function,  $B(z)$ , using Eq. ( 5 )

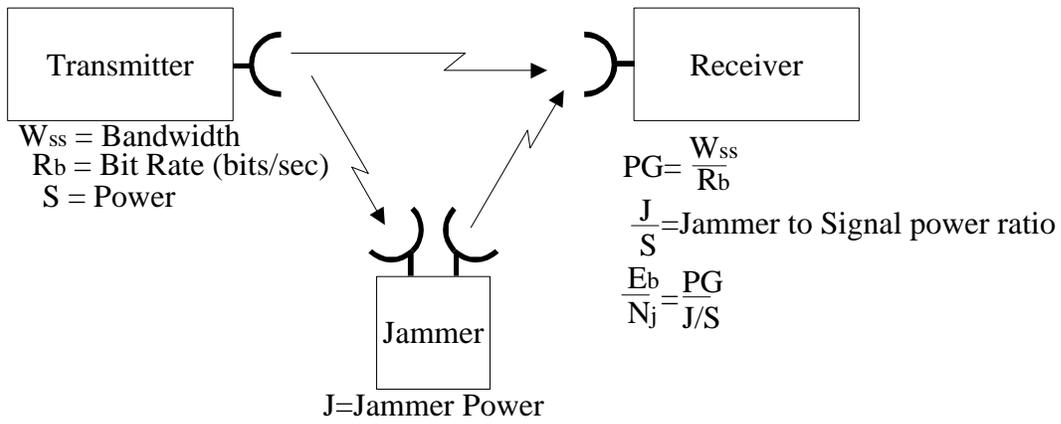


Figure 5. Basic spread spectrum communications system

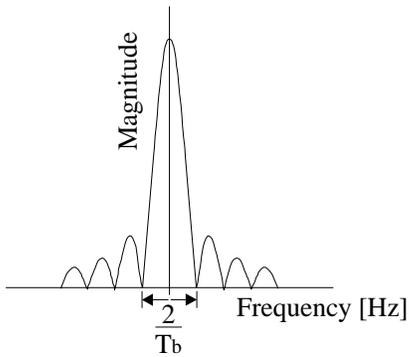


Figure 6. Spectrum of signal BPSK

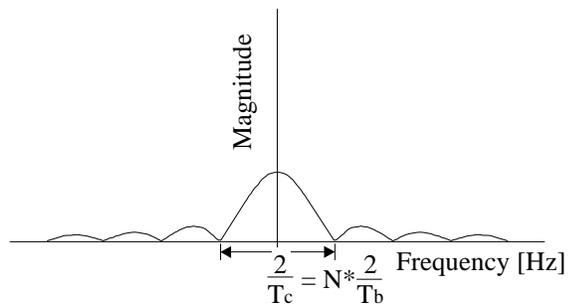


Figure 7. Spectrum of signal BPSK after spreading

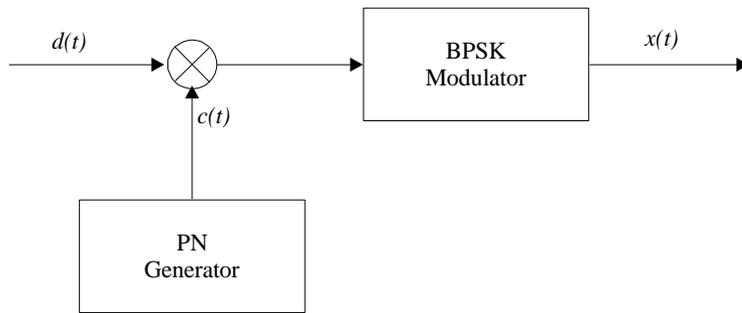


Figure 8. DS/BPSK modulation

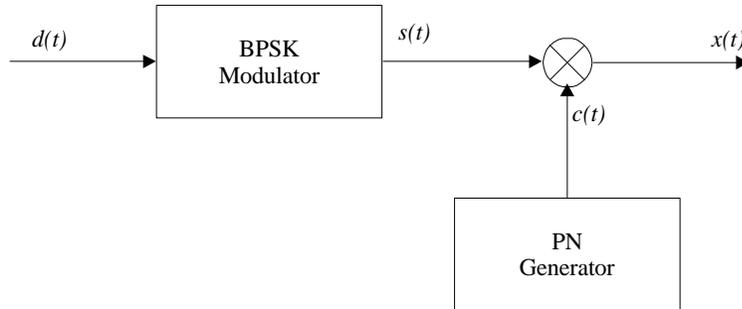


Figure 9. DS/BPSK modified

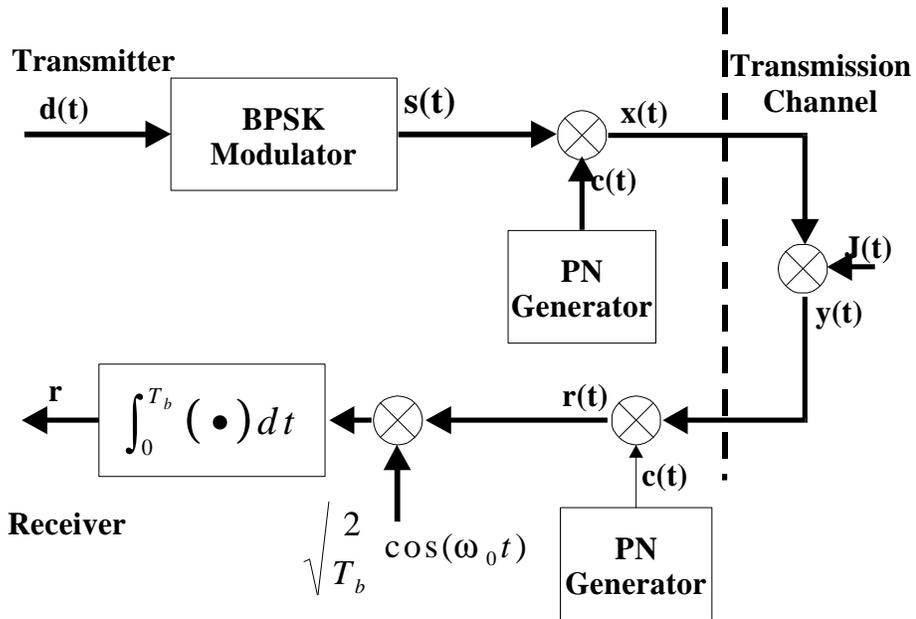


Figure 10. Uncoded DS/BPSK

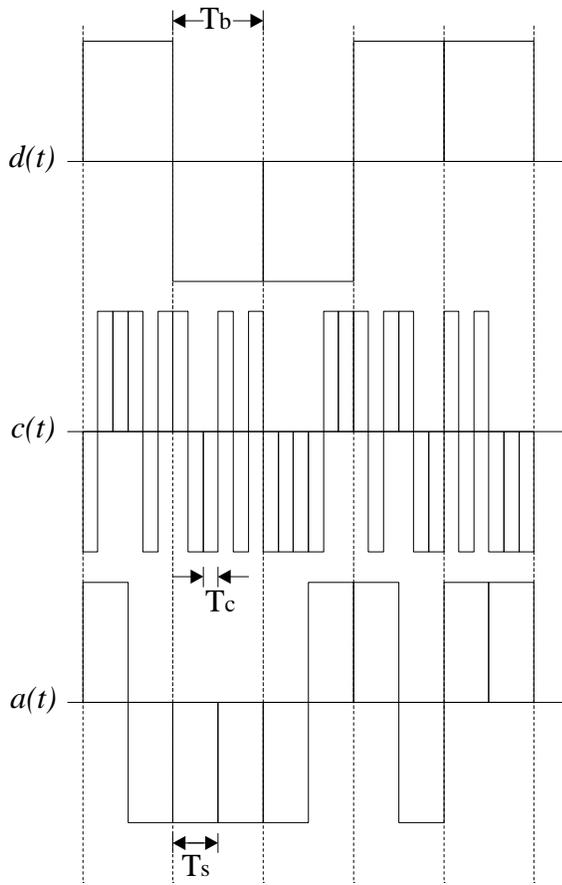


Figure 11. Coded and uncoded signals before spreading

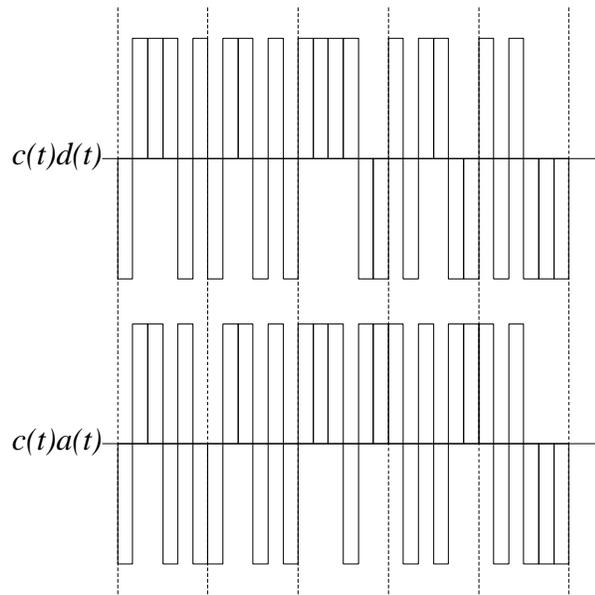


Figure 12. Coded and uncoded signals after spreading

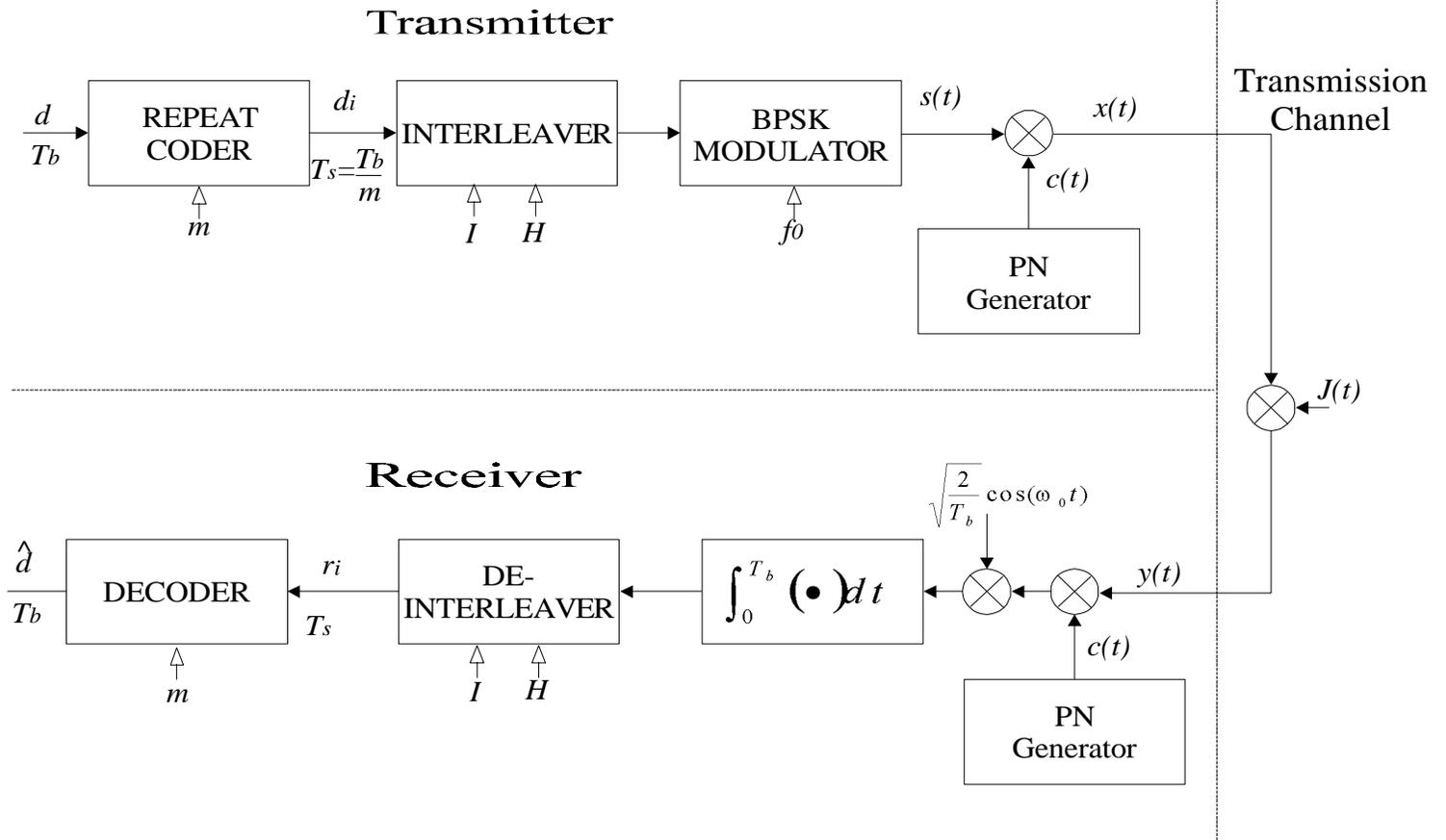


Figure 13. Repeat code DS/BPSK system

$X_1$	$X_{16}$	$X_{31}$	$X_{46}$	$X_{61}$
$X_2$	$X_{17}$	$X_{32}$	$X_{47}$	$X_{62}$
$X_3$	$X_{18}$	$X_{33}$	$X_{48}$	$X_{63}$
$X_4$	$X_{19}$	$X_{34}$	$X_{49}$	$X_{64}$
$X_5$	$X_{20}$	$X_{35}$	$X_{50}$	$X_{65}$
$X_6$	$X_{21}$	$X_{36}$	$X_{51}$	$X_{66}$
$X_7$	$X_{22}$	$X_{37}$	$X_{52}$	$X_{67}$
$X_8$	$X_{23}$	$X_{38}$	$X_{53}$	$X_{68}$
$X_9$	$X_{24}$	$X_{39}$	$X_{54}$	$X_{69}$
$X_{10}$	$X_{25}$	$X_{40}$	$X_{55}$	$X_{70}$
$X_{11}$	$X_{26}$	$X_{41}$	$X_{56}$	$X_{71}$
$X_{12}$	$X_{27}$	$X_{42}$	$X_{57}$	$X_{72}$
$X_{13}$	$X_{28}$	$X_{43}$	$X_{58}$	$X_{73}$
$X_{14}$	$X_{29}$	$X_{44}$	$X_{59}$	$X_{74}$
$X_{15}$	$X_{30}$	$X_{45}$	$X_{60}$	$X_{75}$

Figure 14. Interleaver matrix with  $I=5$  and  $H=15$

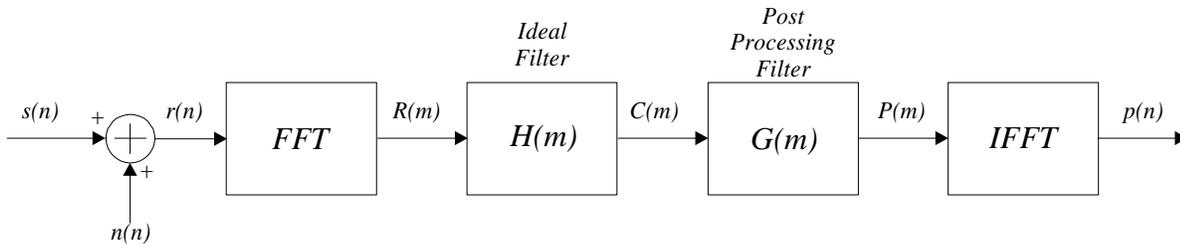


Figure 15. FFT filter assuming additive signal and noise

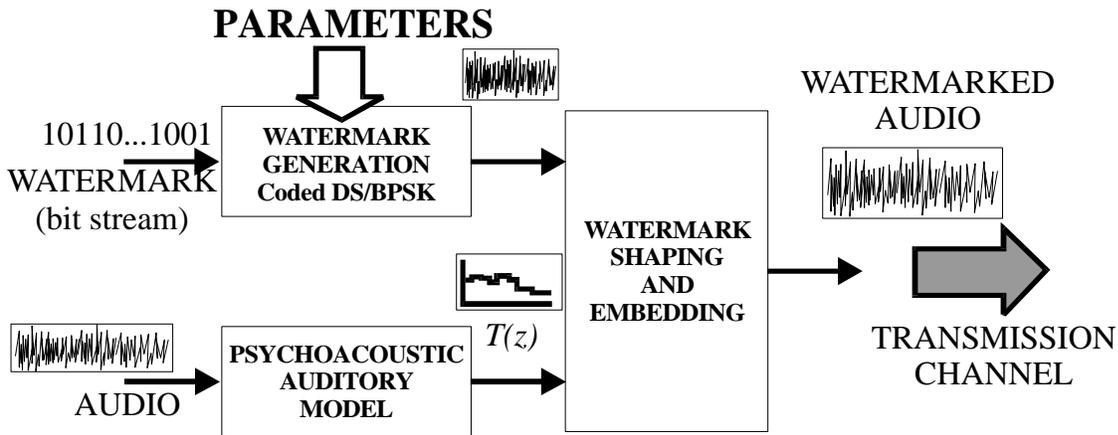


Figure 16 Proposed system (watermark generation and embedding)

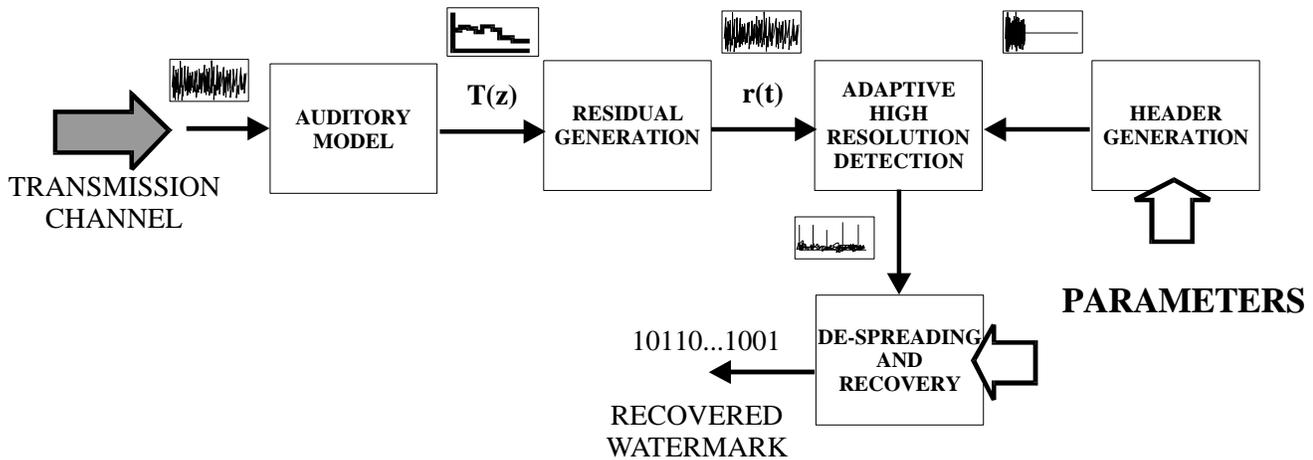


Figure 17 Proposed System (Data recovery)

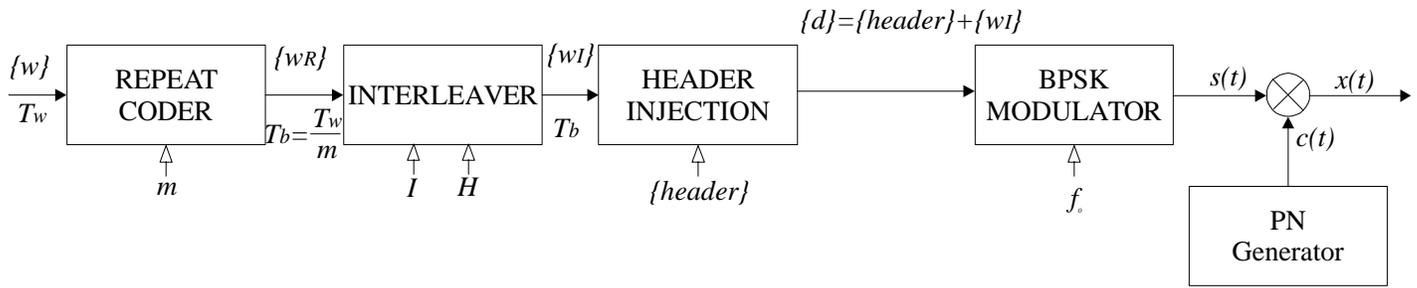


Figure 18. Watermark generation system

1	1	1	-1	1
1	1	-1	-1	1
1	-1	-1	-1	-1
1	-1	-1	1	-1
1	-1	1	1	-1
1	-1	1	1	-1
-1	-1	1	1	-1
-1	-1	1	1	-1
-1	1	1	1	1
1	1	1	1	1

Figure 19. Interleaver matrix

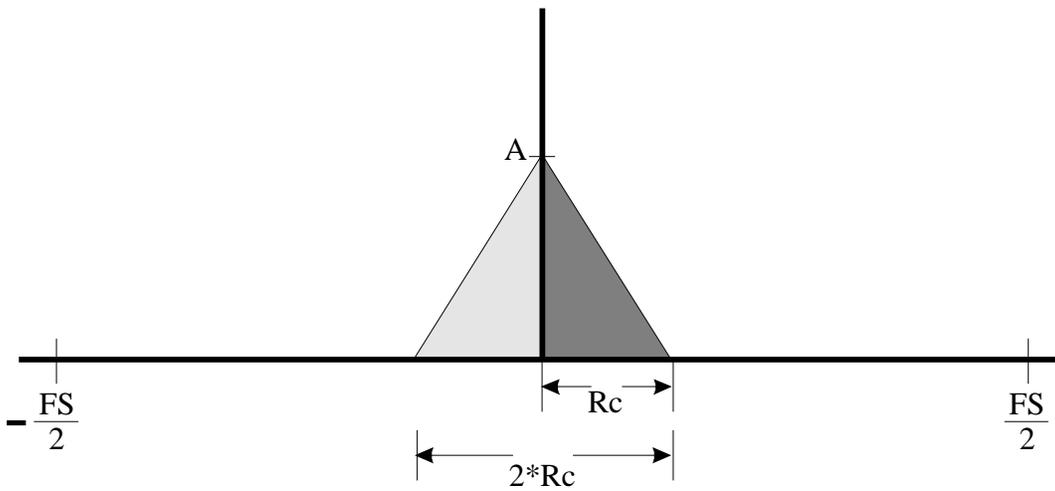


Figure 20. Baseband System Parameters

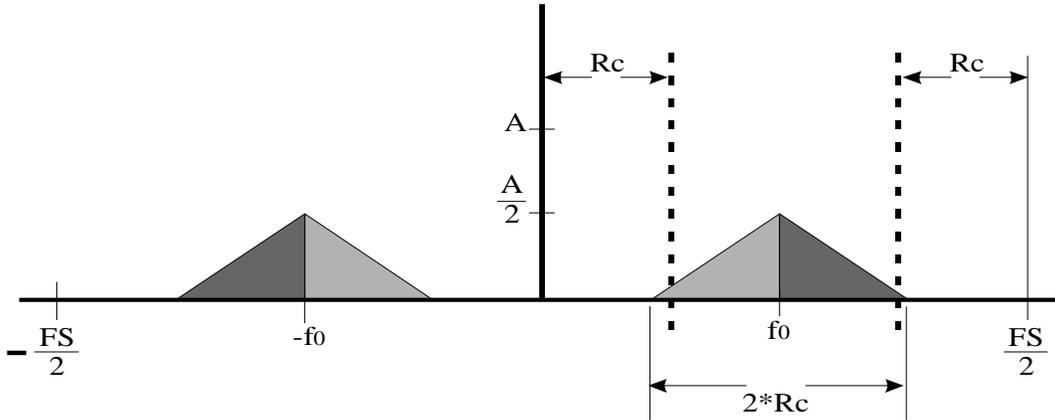


Figure 21. Passband System Parameters (anti-aliasing)

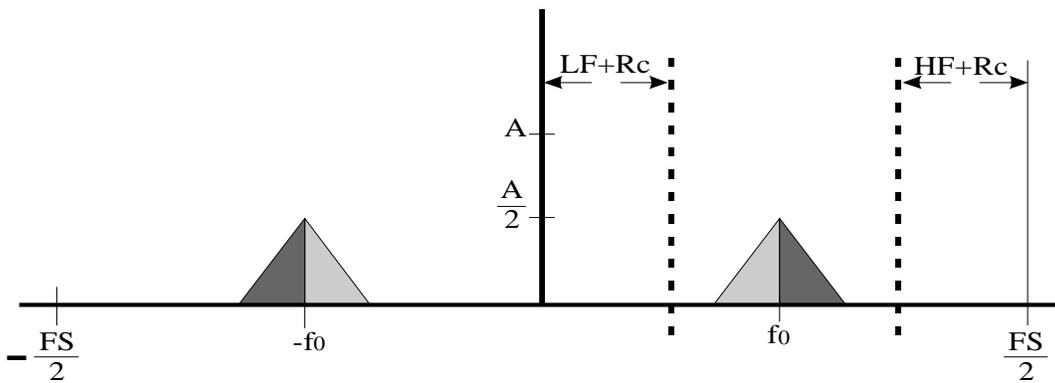


Figure 22. Passband System with Frequency Limits  $LF$  and  $HF$

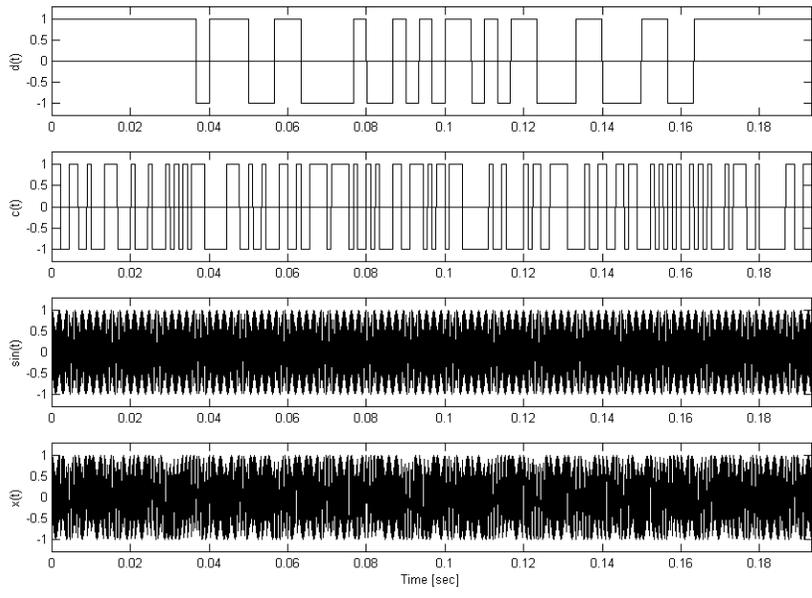
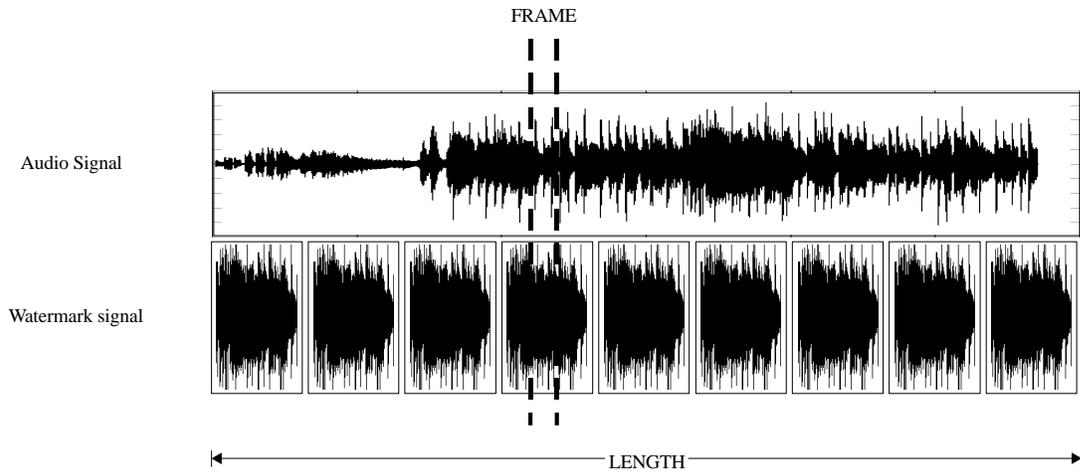
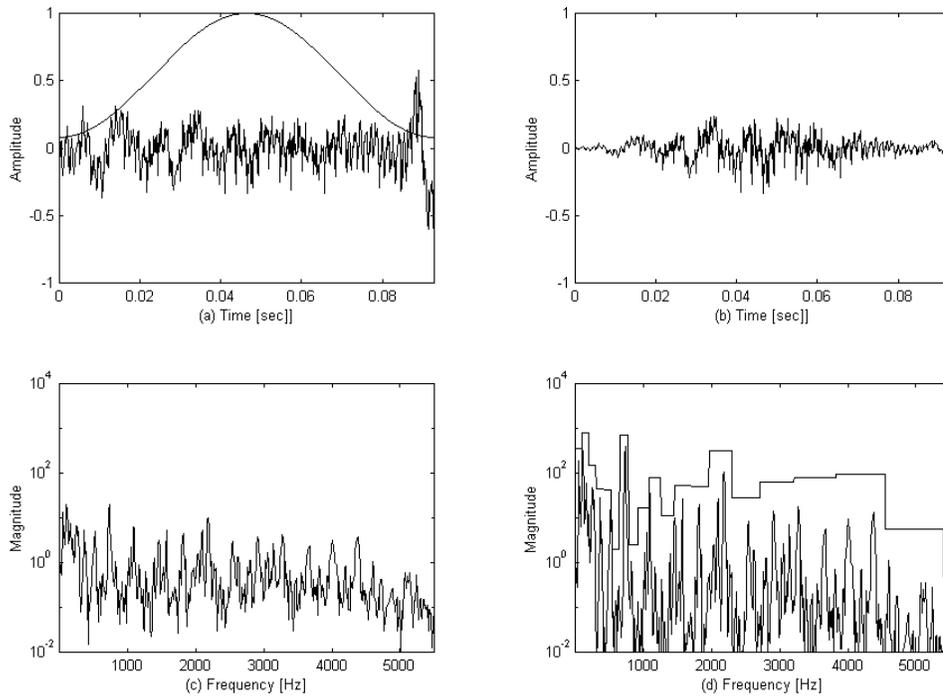


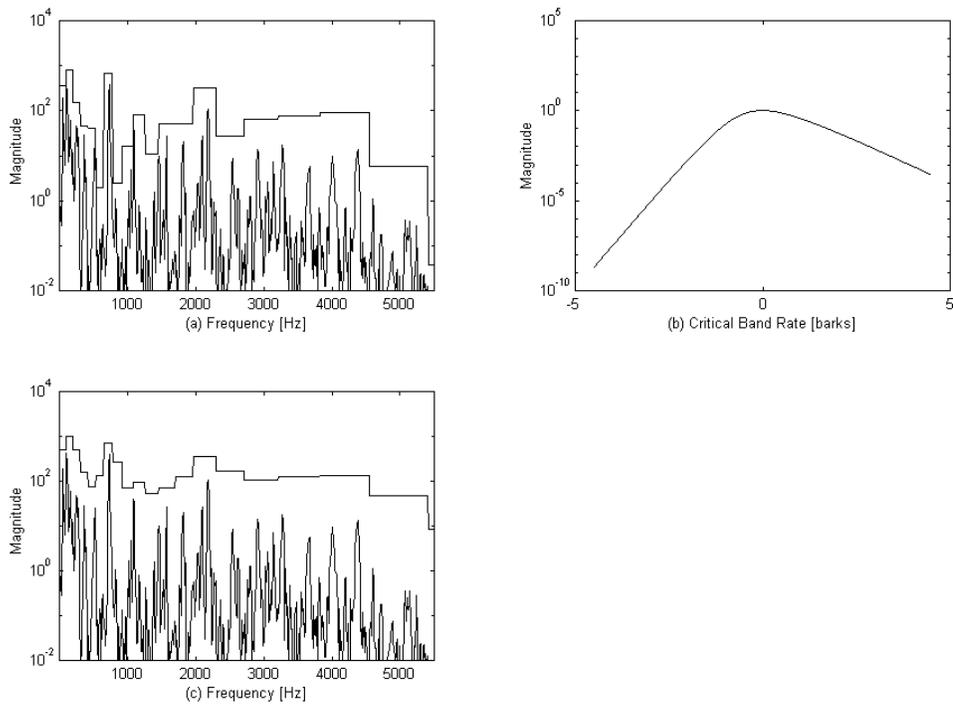
Figure 23 Time domain signals: data bit stream,  $d(t)$ ; PN sequence,  $c(t)$ ; BPSK modulator,  $\sin(t)$ ; and watermark audio signal,  $x(t)$



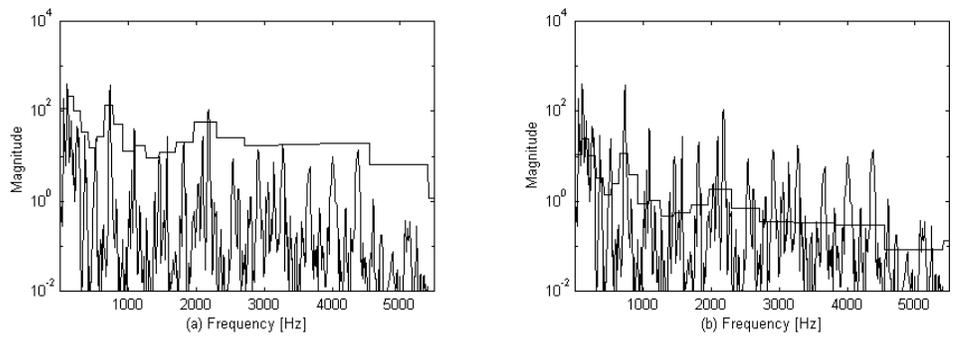
**Figure 24. Frame segmentation and watermark redundancy**



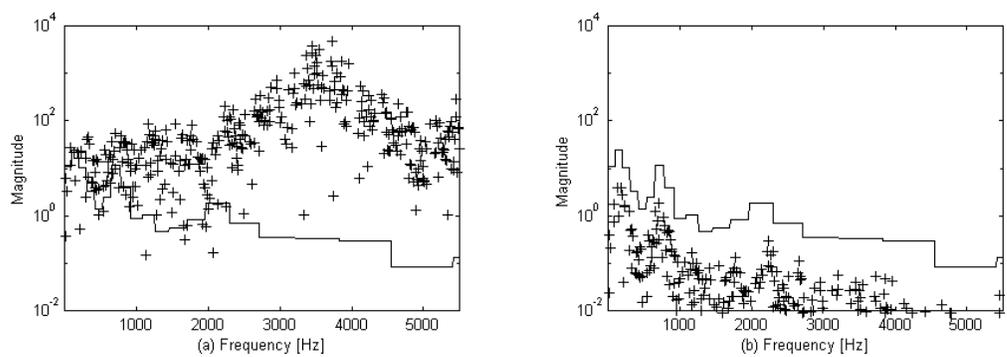
**Figure 25. (a) Audio signal  $s(t)$  and window signal  $w(t)$ , (b) windowed signal  $sw(t)$ , (c) magnitude of frequency representation  $Sw(jW)$ , and (d) power spectrum  $Sp(jW)$  and energy per critical band  $Spz(z)$**



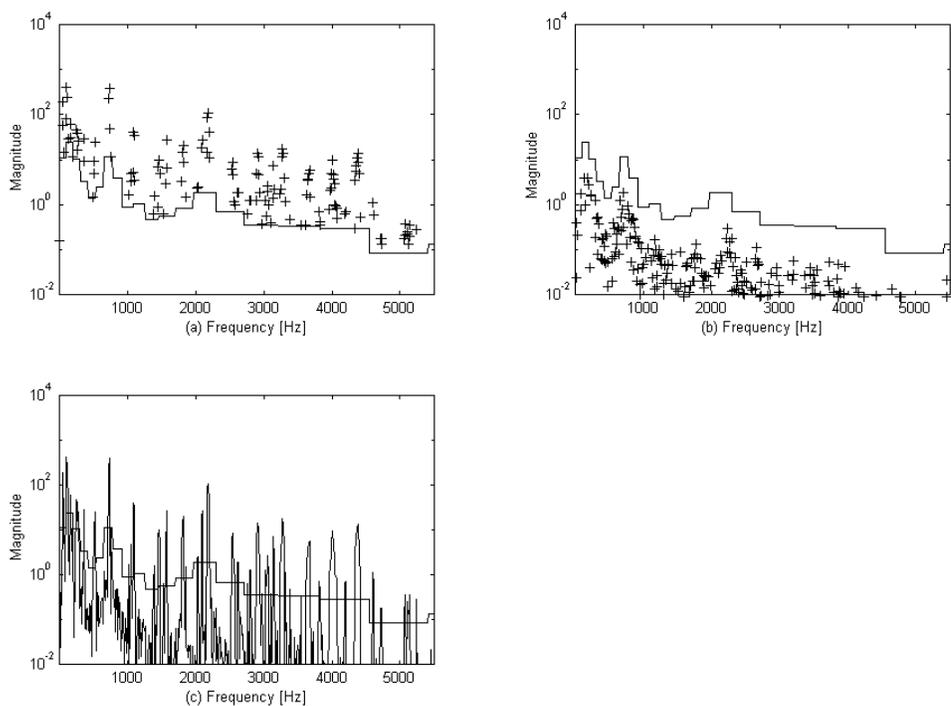
**Figure 26. (a) Energy per critical band  $Spz(z)$ , (b) spreading function  $B(z)$ , and (c) Spread energy per critical band  $Sm(z)$**



**Figure 27. (a) raw masking threshold  $Traw(z)$ , and (b) normalized masking threshold  $Tnorm(z)$**



**Figure 28.** (a)  $X_{wnew}(z)$  before shaping, (b) after shaping with  $A = 0.4$



**Figure 29.** Final masking threshold  $T_{final}(z)$ , and the power spectrum of (a)  $S_{wnew}(j\omega)$ , (b)  $X_{final}(j\omega)$ , and (c)  $OUT(j\omega)$ .

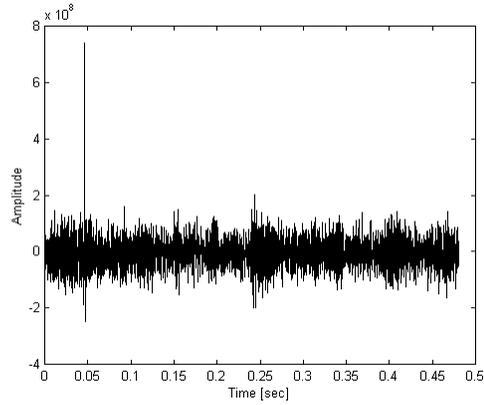


Figure 30. Detection peak in  $det(t)$

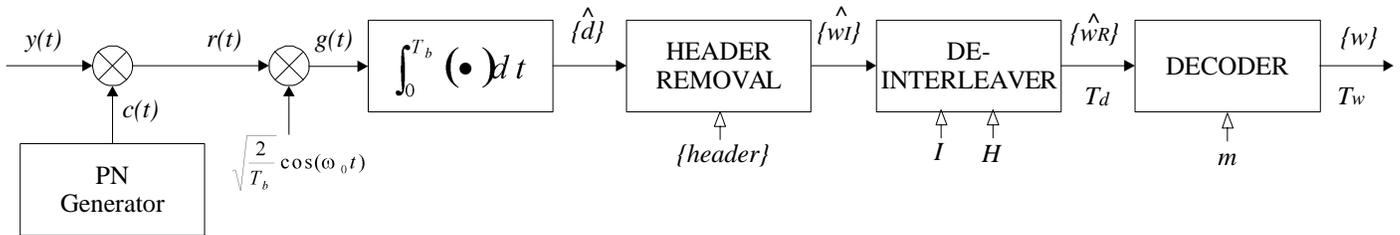


Figure 31. Watermark recovery system

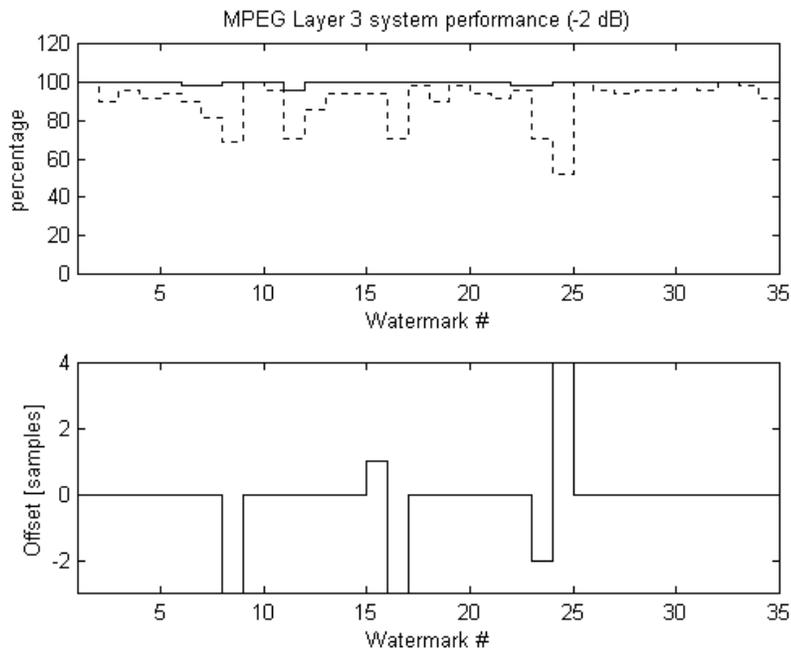
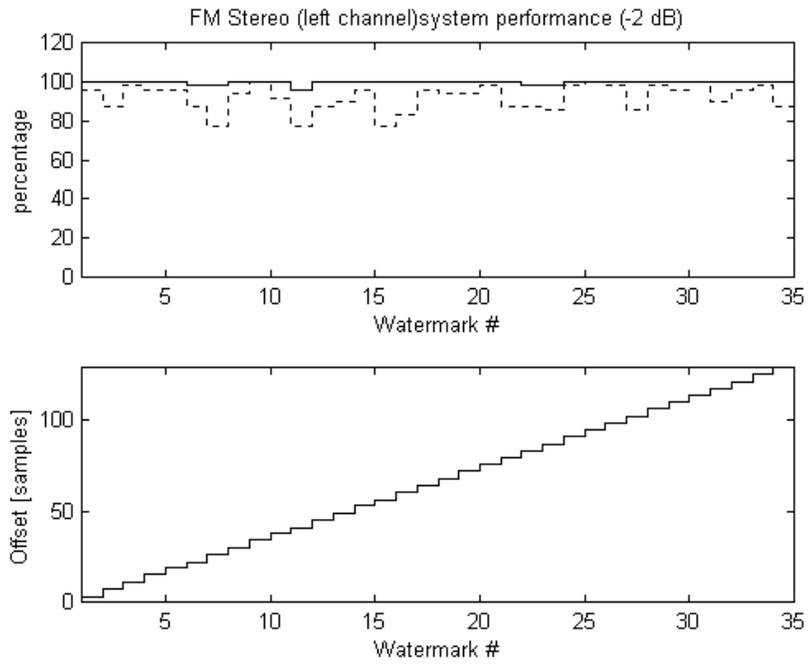


Figure 32 MPEG layer 3 system performance



**Figure 33 FM stereo (left channel) system performance**