WILEY | Hindawi

*Research Article*

# Optimization of a Blind Speech Watermarking Technique against Amplitude Scaling

## Mohammad Ali Nematollahi,[1] Chalee Vorakulpipat,[2] and Hamurabi Gamboa Rosales[3]

[1]*Department of Computer Engineering, Islamic Azad University, Safadasht Branch, Tehran, Iran*
[2]*Cybersecurity Laboratory, National Electronics and Computer Technology Center, 112 Phahonyothin Road, Klong 1,*
 *Klong Luang, Pathumthani 12120, Thailand*
[3]*Department of Electronics Engineering, Universidad Autónoma de Zacatecas, 98000 Zacatecas, DF, Mexico*

Correspondence should be addressed to Chalee Vorakulpipat; chalee.vorakulpipat@nectec.or.th

This paper presents a gain invariant speech watermarking technique based on quantization of the Lp-norm. In this scheme, first, the original speech signal is divided into different frames. Second, each frame is divided into two vectors based on odd and even indices. Third, quantization index modulation (QIM) is used to embed the watermark bits into the ratio of the Lp-norm between the odd and even indices. Finally, the Lagrange optimization technique is applied to minimize the embedding distortion. By applying a statistical analytical approach, the embedding distortion and error probability are estimated. Experimental results not only confirm the accuracy of the driven statistical analytical approach but also prove the robustness of the proposed technique against common signal processing attacks.

## 1. Introduction

Hiding a secret message in an object has a long history, possibly dating back thousands of years. The rapid growth of computer and communication transmissions has inspired the idea of digital data hiding. Digital watermarking, as a major branch of data hiding, has attracted many researchers [1]. The importance of speech watermarking is gradually increasing because of significant speech transmission through insecure communication channels. There are many approaches for speech watermarking, including spread spectrum (SS), auditory masking, patchwork, transformation, and parametric modeling [2]. In the SS approach, a pseudorandom sequence is used to spread the spectrum of the watermark data and add it to the frequency spectrum of the host signal. However, auditory masking uses unimportant perceptual components of the signal to embed the watermark bits. By contrast, the patchwork approach embeds the watermark data by manipulating two sets of the signal to determine the difference between them. The transformation approach embeds the watermark data into the transformation domains, for

example, discrete cosine transform, discrete wavelet transform (DWT), and discrete Fourier transform (DFT). Finally, in the parametric modeling approach, the watermark is embedded by modifying the coefficients of the autoregressive (AR) model.

In addition to speech watermarking approaches, four main embedding strategies are widely applied for watermarking: least significant bit (LSB) replacement, quantization, addition, and multiplication. Among these strategies, quantization has attracted much attention because of blindness, robustness, controlled distortion, and payload. For this purpose, a set of quantizers that are associated with various watermark data are used. However, the quantization strategy suffers from amplitude scaling. To rectify this problem, rational dither modulation (RDM) [3] was proposed to enhance the robustness of quantization index modulation (QIM) [4, 5]; however, it degraded the imperceptibility of the watermarked signal. Hence, hyperbolic RDM [6] was proposed to improve the robustness against power law and gain attacks. Another attempt was made by embedding a watermark into the angle of the signal, known as angle QIM

(AQIM) [7]. However, this technique was very sensitive to additive white Gaussian noise (AWGN). In [8], normalized cross-correlation between the original signal and a random sequence was quantized based on dither modulation (known as NC-DM) to embed the watermark data. However, applying the random sequence degraded the security of this technique. Lastly, other efforts, such as projection quantization [9], logarithmic quantization index modulation (LQIM) [10], and Lp-norm QIM [11], have been studied for a gain invariant image watermarking technique.

This paper attempts to mitigate the limitations of previous research by quantizing the ratio between the Lp-norms of even and odd indices. After quantization, the Lagrange optimization method is applied to compute the best watermarked sample that minimizes the embedding distortion and improves imperceptibility. By assuming Laplacian and Gaussian distributions for the speech and noise signals, respectively, the embedding distortion and error probability are driven analytically and validated by performing a simulation. Moreover, experimental results show that the proposed speech watermarking technique outperforms state-of-the-art watermarking techniques.

Generally, speech watermarking should preserve the identity of the speaker, which is important for certain security applications [12, 13]. To preserve speaker-specific information, some investigations have been conducted to embed the watermark into special frequency subbands that have less speaker-specific information [5, 14, 15]. Further discussion can be found in [16].

The remainder of this paper is organized as follows. In Section 2, the proposed model for the speech watermarking technique is presented. Additionally, the watermark embedding and extraction processes are described. The performance of the developed watermarking technique is analytically studied in Section 3 and validated by performing a simulation in Section 4. The experimental results are explained in Section 5. Finally, the conclusion and future work are discussed.

## 2. Proposed Speech Watermarking Technique

In this section, a blind speech watermarking technique is developed based on quantization of the Lp-norm ratio between two blocks of even and odd indices. Assume that $S$ represents an original speech signal that consists of $N$ samples. Two subsets $X$ and $Y$ are formed with respect to even and odd indexed terms, respectively, so that both $X$ and $Y$ have approximately the same energy that causes less embedding distortion. Moreover, synchronization between the transmitter and receiver is most efficient in this case. Figure 1 shows the formation of the subsequences of $X$ and $Y$ from the odd and even indices of the original signal, respectively.

Then, the Lp-norm of both subsequences $X$ and $Y$ are computed, respectively, as follows:

$$L_X = \sqrt[P]{\frac{2}{N} \times \sum_{i=1}^{N} |S_{2i}|^P}, \tag{1}$$

$$L_Y = \sqrt[P]{\frac{2}{N} \times \sum_{i=1}^{N} |S_{2i-1}|^P}. \tag{2}$$

The ratio ($Z$) between $L_X$ and $L_Y$, given as

$$Z = \frac{L_X}{L_Y}, \tag{3}$$

is quantized to embed the watermark bit. Although embedding the watermark into the ratio of Lp-norms can provide high robustness against various attacks, imperceptibility can be seriously degraded.

To resolve this limitation, the variation between the original ratio ($Z$) and quantized ratio ($Z^Q$) should be minimized. Therefore, the Lagrange optimization method is used to minimize this variation; that is, the Lagrange optimization method decreases the embedding distortion after quantization to improve the imperceptibility of the watermarked speech signal. As a result, the Lagrange optimization problems can be formulated as follows:

$$\text{Minimize:} \quad J(X) = \sum_{j=1}^{N/2} \left( X_j^Q - X_j \right)^P,$$

$$\text{Subject to:} \quad C(Y) = \sqrt[P]{\frac{2}{N} \sum_{j=1}^{N/2} \left( X_j^Q \right)^P} - Z^Q \times L_Y = 0. \tag{4}$$

To solve this optimization problem, the Lagrange method should estimate the optimized values of the equation system as follows:

$$\nabla J(X) - \lambda \nabla C(Y) = 0. \tag{5}$$

These optimized values are simply computed by solving the following:

$$X_j^{Q,\text{opt}} = \lambda_{\text{opt}} \times X_j, \tag{6}$$

$$\lambda_{\text{opt}} = \frac{Z^Q \times L_Y}{L_X}. \tag{7}$$

*2.1. Speech Watermarking Algorithm.* The details of the proposed embedding and extraction processes are described in the following algorithms:

*Embedding Process*

(a) Segment the input speech signal ($S$) into different frames ($S_i$) with size $N$.

(b) Form two subsequences $X$ and $Y$, each of length $N/2$, based on the even and odd indices of $S_i$, respectively.

(c) Compute the Lp-norms $L_X$ and $L_Y$ of both the $X$ and $Y$ subsequences, respectively, based on (1) and (2), respectively.

(d) Apply the QIM technique to embed the watermark bit into the ratio between the Lp-norms of $X$ and $Y$ ($Z = L_X/L_Y$) as follows:

$$Z^Q = \left\lfloor \frac{Z + W_i \times \Delta}{2\Delta} \right\rfloor \times 2\Delta + W_i \times \Delta, \quad W_i \in \{0, 1\}, \tag{8}$$
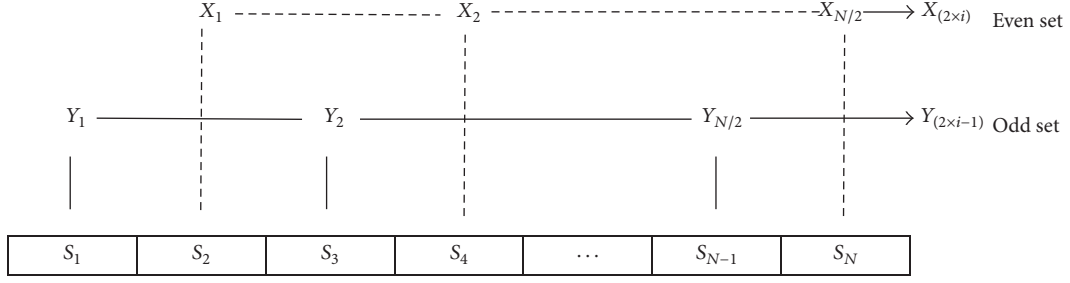
FIGURE 1: Formation of two odd and even subsets from the original speech signal.

where $\Delta$ represents the quantization steps, $W_i$ is the watermark bit, and $Z^Q$ is the modified ratio of the Lp-norms between $X$ and $Y$. Choosing large quantization steps increases the robustness but results in less imperceptibility and vice versa.

(e) Apply the Lagrange method to optimize the values of $X_{2i}^Q$.

(f) Reposition the even and odd subsequences based on $X_{2i}^Q$ and $Y$, respectively.

(g) Rearrange the watermarked speech signal based on the modified frames ($\widehat{S}_i$).

Figure 2 shows the block diagram of the proposed embedding process.

*Extraction Process*

(a) Segment the input watermarked speech signal ($\widehat{S}$) into different frames ($\widehat{S}_i$) with size $N$.

(b) Form two subsequences $\widehat{X}$ and $\widehat{Y}$, each of length $N/2$, based on the even and odd indices of $S_i$, respectively.

(c) Compute the Lp-norms $\widehat{L_X}$ and $\widehat{L_Y}$ of both $\widehat{X}$ and $\widehat{Y}$ subsequences, respectively, based on (1) and (2), respectively.

(d) Extract the $k$th binary watermark data from the $k$th frame of the watermarked speech signal by selecting the minimum Euclidean distance (nearest quantization step) from the ratio of $\widehat{Z}_k = \widehat{L_X}/\widehat{L_Y}$ as follows:

$$\widehat{w}_k = \min\left(\sqrt{\widehat{Z}_k^2 + Q_0\left(\widehat{Z}_k\right)^2}, \sqrt{\widehat{Z}_k^2 + Q_1\left(\widehat{Z}_k\right)^2}\right), \quad (9)$$

where $Q_{b_k}$ is the quantization function while meeting the requirements of watermark bits $b_k = \{0, 1\}$.

Figure 3 shows the block diagram of the proposed extraction process.

## 3. Statistical Analysis of the Proposed Technique

Generally, Laplacian distribution is the best distribution approach for modeling speech signals within the frame range of 5–50 ms [17, 18]. Laplacian distribution is expressed as

$$f(x) = \frac{b}{2}e^{(-b|x-\mu|)}, \quad b = \frac{L}{\sum_{i=1}^{L}|x_i - \mu|}, \quad (10)$$

where $L$ is the sample size and $\mu$ is the mean of the random variables. If the subsequences of $X$ and $Y$ are considered as independent, identically distributed (i.i.d) variables, then the distribution of each of them can be assumed to be Laplacian distributions $X = \mathcal{CL}(\mu_X, 2b_X^2)$ and $Y = \mathcal{CL}(\mu_Y, 2b_Y^2)$, respectively. Based on (3), the ratio ($Z$) between $X$ and $Y$ should be computed. However, the ratio between two Laplacian distributions cannot be computed exactly because the mean and variance are not actually finite in either the Gaussian or Laplacian case. The problem arises because the denominator has nonzero density in the neighborhood of zero. If the denominator is bounded away from zero (immediately it no longer has the ratio of two Laplacian distributions or two normals), then a Taylor expansion should converge to estimate the ratio between two Laplacian distributions. According to Appendix A, the parameters of the ratio can be derived as follows:

$$\mu_Z = \frac{\mu_X}{\mu_Y}\left(1 + \frac{2b_Y^2}{\mu_Y^2}\right), \quad (11)$$

$$\sigma_z^2 = \frac{\mu_X^2}{\mu_Y^2}\left(\frac{2b_Y^2}{\mu_Y^2} - \frac{4b_Y^4}{\mu_Y^4}\right) + \frac{2b_X^2}{\mu_Y^2}. \quad (12)$$

To estimate the embedding distortion, quantization noise ($\Delta$) should be considered between the original and watermarked speech signals as follows:

$$S_i - \widehat{S}_i = \left(\widehat{X}_{2i} - \widehat{Y}_{2i-1}\right) - \left(X_{2i} - Y_{2i-1}\right). \quad (13)$$

As in (4) to (6), $\widehat{Y}_{2i-1} = Y_{2i-1}$; thus, (12) can be expressed as

$$S_i - \widehat{S}_i = \lambda_{\text{opt}} \times X_{2i} - X_{2i} = X_{2i} \times \left(\lambda_{\text{opt}} - 1\right). \quad (14)$$

If $Z_i^Q = (L_X + \varepsilon)/L_Y$, then $\lambda_{\text{opt}}$ can be expressed as

$$\lambda_{\text{opt}} = \left(\frac{L_X + \varepsilon}{L_Y}\right) \times \frac{L_Y}{L_X} = \left(1 + \frac{\varepsilon}{L_X}\right). \quad (15)$$
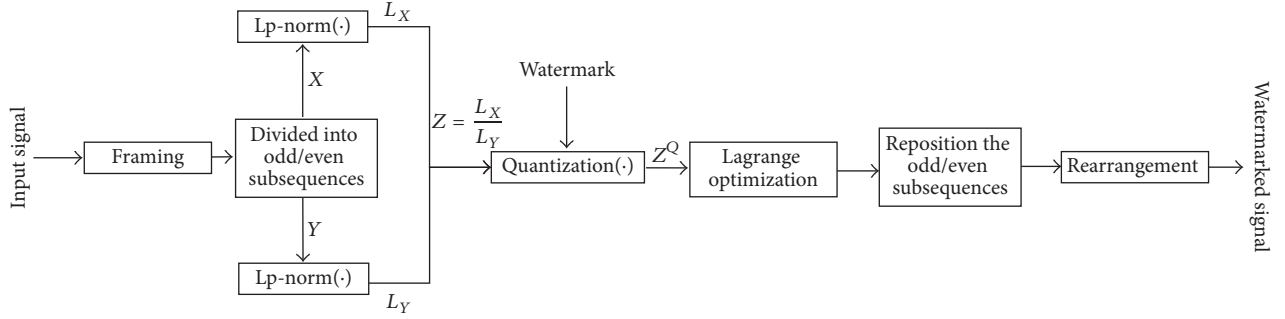
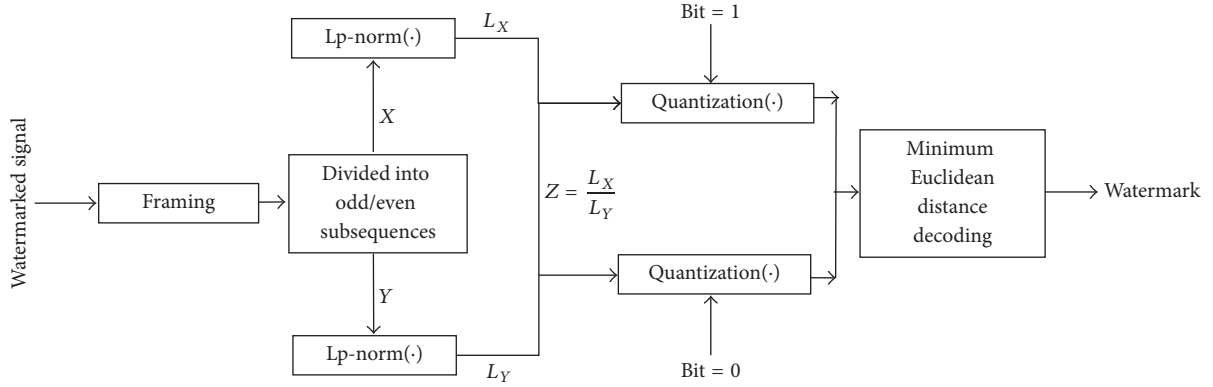FIGURE 2: Block diagram of the proposed embedding process.



FIGURE 3: Block diagram of the proposed extraction process.

Thus, (13) can be approximately estimated by

$$S_i - \widehat{S_i} \approx \left(1 + \frac{\varepsilon}{L_X}\right) \times X_{2i}, \tag{16}$$

Therefore, the expected values of (13) can be estimated as

$$E\left[\left\|\widehat{S} - S\right\|^2\right] \cong E\left(\left(1 + \frac{\varepsilon}{L_X}\right)^2\right) E\left(X_{2i}^2\right)$$

$$= \left[E\left(\frac{1}{L_X^2}\right)E\left((\varepsilon)^2\right) + E\left(\varepsilon\right)E\left(\frac{2}{L_X}\right) + 1\right]E\left(X_{2i}^2\right). \tag{17}$$

If quantization noise ($\varepsilon$) is considered as a uniform distribution in $[-\Delta/2, \Delta/2]$ then $E(\varepsilon) = 0$ and $E((\varepsilon)^2) = \Delta^2/48$. Additionally, as the mean value of the speech signal is considered to be zero, then the zero mean Laplacian distribution is used to model the speech signal as $E(X_{2i}) = 0$. As a result, $(X_i^2) = 2b_{X_i}^2$. To model $E(1/L_X^2)$, the absolute moment of the Laplacian distribution should be estimated using Appendix B as follows:

$$E\left(|X|^P\right) = \left(\frac{e^{\mu/b}b^n}{2}\right)\left[(-1)^n \cdot I_n + n!\right], \tag{18}$$

where $I_n = nI_{n-1}$ and $I_n = \int_{-\infty}^{0} t^n e^{-t} dt$. Thus, we can derive the mean and variance for the $P$th absolute moment of the Laplacian distribution as

$$L_X^P = \sum_{j=1}^{N/2} \left|X_j\right|^P \sim \mathcal{CL}\left(\mu_{XP}, \frac{2\left(\mu_{X(2P)} - \mu_{XP}^2\right)}{N}\right). \tag{19}$$

Now, based on (1) and (19), we can compute $E(1/L_X^2) = E(1/(\sqrt[p]{L_X^P})^2) = \mu_{X(2/P)}$. Therefore, the signal-to-watermark ratio (SWR) can be estimated as

$$\text{SWR} = \frac{E\left[\|S\|^2\right]}{E\left[\left\|\widehat{S} - S\right\|^2\right]}$$

$$\cong \frac{2b_X^2 + 2b_Y^2}{\left((\Delta^2/48) \times \mu_{X(2/P)} + 1\right) \times (2b_X^2)}. \tag{20}$$

Because both $X$ and $Y$ sets have been selected from the neighboring samples, it can be assumed that $2b_X^2 \cong 2b_Y^2$. As a result, (20) can be expressed based on the quantization step as

$$\Delta = \sqrt{\frac{-2\left(1 + 12 \times \text{SWR}\right)}{\text{SWR} \times \mu_{X(2/P)}}}. \tag{21}$$

To model the error probability, it is assumed that the watermarked speech signal passes through an AWGN channel with zero mean Gaussian noise $\mathcal{N}(0, \sigma_n^2)$. Therefore, (3) must be rewritten as

$$\widehat{Z} = \frac{\sum_{j=1}^{N/2} \left|X_j + N_{X_j}\right|^P}{\sum_{j=1}^{N/2} \left|Y_j + N_{Y_j}\right|^P}, \tag{22}$$

where $N_{Y_j}$ and $N_{X_j}$ correspond to the odd and even components of the AWGN, respectively. Because the term

$\sum_{j=1}^{N/2} |X_j|^P$ is a known parameter, it is not possible to estimate $\widehat{Z}$ using a chi-square with $N$ degrees of freedom, $\varkappa^2(N)$. To compute the distribution of $\widehat{Z}$, it should be decomposed and estimated as

$$\widehat{Z} = \frac{\sum_{j=1}^{N/2} \left( |X_j|^P + P |X_j|^{P-1} N_{X_j} + (P(P-1)/2) |X_j|^{P-2} N_{X_j}^2 + \cdots + N_{X_j}^P \right)}{\sum_{j=1}^{N/2} \left( |Y_j|^P + P |Y_j|^{P-1} N_{Y_j} + (P(P-1)/2) |Y_j|^{P-2} N_{Y_j}^2 + \cdots + N_{Y_j}^P \right)}. \tag{23}$$

Equation (23) can be expressed as

$$\widehat{Z} \approx \text{Original } Z + \overbrace{\gamma_1 + \gamma_2 + \gamma_3}^{\text{Noise}}, \tag{24}$$

where each part of $\widehat{Z}$ is estimated as follows:

$$\text{Original } Z = \frac{\sum_{j=1}^{N/2} |X_j|^P}{\sum_{j=1}^{N/2} |Y_j|^P},$$

$$\gamma_1 = \frac{\sum_{j=1}^{N/2} P |X_j|^{P-1} N_{X_j}}{\sum_{j=1}^{N/2} |Y_j|^P},$$

$$\gamma_2 = -\frac{\sum_{j=1}^{N/2} |X_j|^P}{\sum_{j=1}^{N/2} |Y_j|^P} \times \frac{\sum_{j=1}^{N/2} P |Y_j|^{P-1} N_{Y_j}}{\sum_{j=1}^{N/2} |Y_j|^P}, \tag{25}$$

$$\gamma_3 = \frac{\sum_{j=1}^{N/2} P |X_j|^{P-1} N_{X_j}}{\sum_{j=1}^{N/2} |Y_j|^P}$$
$$\times \frac{\sum_{j=1}^{N/2} P |Y_j|^{P-1} N_{Y_j}}{\sum_{j=1}^{N/2} P |Y_j|^{P-1} N_{Y_j}}.$$

To estimate the probability of error, the noise term can be analyzed because it makes the original $Z$ into a wrong region. Therefore, the distribution of each term of (24) can be estimated by the central limit theorem (CLT) because of the large number of samples in each block. Regardless of the type of original speech signal distribution and because of the independence between the signal and noise samples, the mean and variance of the noise can be computed as

$$\mu_{\text{Noise}} = \mu_{\gamma 1} + \mu_{\gamma 2} + \mu_{\gamma 3},$$
$$\sigma_{\text{Noise}}^2 = \sigma_{\gamma 1}^2 + \sigma_{\gamma 2}^2 + \sigma_{\gamma 3}^2. \tag{26}$$

By assuming equal probabilities for both zero and one bit of the watermark data, the probability of error for a fixed quantization step ($\Delta$) can be estimated as

$$P_e = \sum_{i=1}^{\infty} \frac{1}{2} \Pr \left\{ T_{(i-1)/2} < Z^P < T_{(i+1)/2} \right\}$$
$$\times \sum_{j=-\lfloor i/2 \rfloor}^{\infty} \Pr \left\{ V_{2j+i} < \widehat{Z}^P < V_{2j+i+1} \right\}. \tag{27}$$

A close-form solution for (27) is computed as

$$P_e = \sum_{i=1}^{\infty} \left( Q \left( \frac{T_{(i-1)/2 - \mu_Z}^P}{\sigma_Z} \right) - Q \left( \frac{T_{(i+1)/2 - \mu_Z}^P}{\sigma_Z} \right) \right)$$
$$\times \sum_{j=-\lfloor i/2 \rfloor}^{\infty} \left( Q \left( \frac{V_{(i+2j)/2 - \mu_{\widehat{Z}^P}}^P}{\sigma_{\widehat{Z}^P}} \right) \right. \tag{28}$$
$$\left. - Q \left( \frac{V_{(i+2j+1)/2 - \mu_{\widehat{Z}^P}}^P}{\sigma_{\widehat{Z}^P}} \right) \right),$$

where $Q(\cdot)$ is the complementary error function defined as $Q(x) = (1/\sqrt{2\pi}) \int_x^{\infty} e^{-u^2/2} du$, $T_i = i\Delta$, $V_i = (T_{i/2} + T_{(i+1)/2})/2$, and $\mu_Z$ and $\sigma_Z$ can be computed as in (11) and (12), respectively.

## 4. Discussion on the Experimental Results

To validate the performance of the developed watermarking technique, a simulation was performed on the TIMIT database to verify the robustness, imperceptibility, and capacity of the technique. The TIMIT database included 630 speakers (438 males and 192 females) with sampling frequency 16 KHz [19]. Each speaker pronounced 10 sentences, which are contained in 6,300 sentences. For the experimental results, the average results of 630 speech signals with duration 1 s to 3 s from 630 speakers were used.

Figure 4 shows the bit error rate (BER) with respect to different $P$ for various frame lengths under Watermark to Noise Ratio (WNR) = 40 dB. In this figure, each curve is plotted separately in order to appear the changes. As can be observed, the frame size was positively correlated with the BER. Whenever the frame size decreased, the BER increased. Additionally, it seems that $P$ was not highly correlated with the BER for $P$ values greater than two. Only a small fluctuation can be observed for the BER when $P$ changed.

Figure 5 shows the BER with respect to different $P$ for various quantization steps. As expected, whenever the quantization step increased, the BER decreased. Furthermore, the variation of $P$ did not seriously change the BER. It must be mentioned that because of perfect watermark detection under clean conditions, a small AWGN was induced on the watermarked signals for the experiments shown in Figures 4 and 5.

Figure 6 shows the variation of the signal-to-noise ratio (SNR) with respect to different $P$ values for different frame lengths. There was not a significant difference in the SNR when the frame size increased. As can be observed, whenever
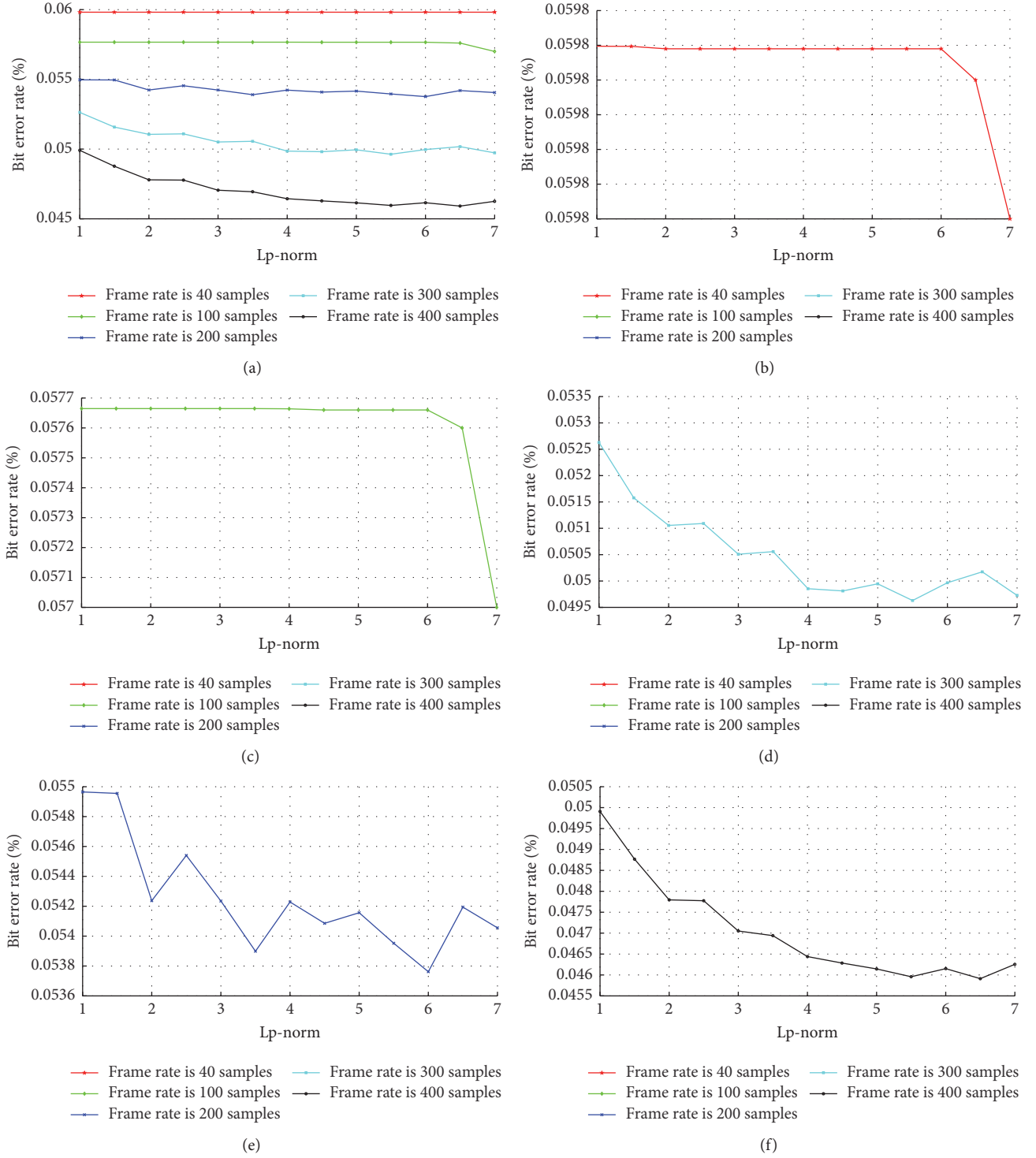
(a)



(b)



(c)



(d)



(e)



(f)

FIGURE 4: (a) BER versus Lp-norms for different frame lengths under WNR = 40 dB (b–f) each curve separately.

the frame size increased, the energy level between the two sets of $L_X$ and $L_Y$ increased. Consequently, the ratio between them increased, which caused a lower SNR. Additionally, it seems that changing $P$ was not highly correlated with the SNR for different frame lengths.

Figure 7 illustrates different SNRs with respect to different $P$ for various quantization steps. As observed, $P$ did not highly affect the SNR. However, the quantization step highly affected the SNR. As the quantization step increased, the SNR decreased.
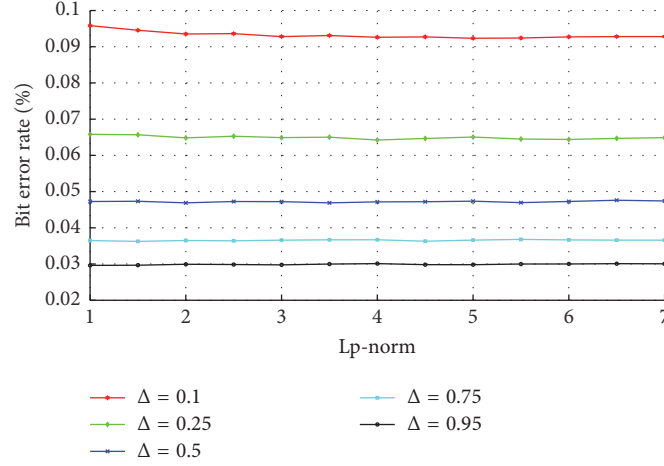
FIGURE 5: BER versus Lp-norms for different quantization steps under WNR = 40 dB.
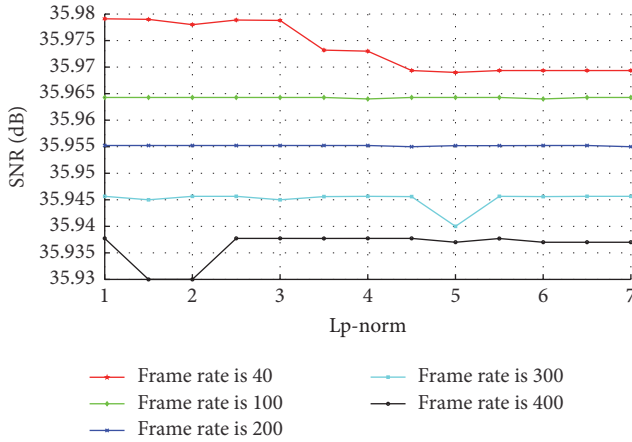


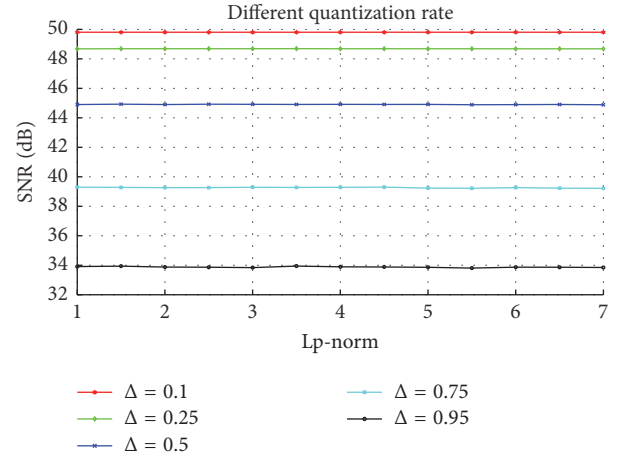FIGURE 6: SNR versus Lp-norms for different frame lengths.



FIGURE 7: SNR versus Lp-norms for various quantization steps.

To compute the payload of the proposed watermark, a memoryless binary symmetric channel (BSC) ($C_{BSC}$) defined as

$$C_{BSC} = R \times [1 + H(P_e)], \tag{29}$$

where

$$H(P_e) = P_e \times \log_2^{(P_e)} + (1 - P_e) \times \log_2^{(1-P_e)}, \tag{30}$$

was applied to estimate the capacity of the channel with bitrate ($R$) for error-free watermark transmission [20].

Because the sampling rate of the TIMIT was 16 KHz, $R$ was assumed to be 64 Kbps (8 KHz for speech bandwidth × 8 bits per sample = 64 Kbps) for a telephony channel and $P_e$ was assumed to be equal to the BER in the watermark detection process. Figure 8 shows the amount of the BSC for different WNRs for various quantization steps. As observed, the capacity increased whenever the WNR increased. This is because the watermark was extracted with a minimum BER when the WNR increased. Moreover, it can be inferred that the amount of the BSC increased while the quantization step

increased because the watermark was embedded with high intensity when the quantization step increased. As observed, the BSC capacity for fewer quantization steps ($\Delta \leq 0.25$) was approximately zero under a high noisy channel.

Figure 9 shows the variation of the BSC capacity with respect to different WNRs for different frame lengths. As observed, it seems that, under serious noise, the frame size was not a significant factor for the BSC capacity. Despite this, the frame size was likely to be important whenever the WNR increased. Thus, for a large WNR, it is obvious that whenever the frame size increased, the BER in the watermark detection process decreased, which caused an improvement in the BSC capacity.

To demonstrate the efficiency and performance of the proposed speech watermarking technique, the robustness, capacity, and inaudibility of the proposed technique must be compared with other state-of-the-art speech watermarking techniques.

Table 1 describes the benchmark for simulating the results for the robustness test. Many of these attacks are based on the StirMark Benchmark for Audio (SMBA) [24].

Table 1: Benchmark for speech watermarking.

| Attack type | Attack name | Description | Parameter(s) | Default value(s) | |
|---|---|---|---|---|---|
| Additive Noise | AddBrumm | It adds buzz or low frequency sinus tone to the watermarked signal to simulate the impact of a power supply | ⟨STRENGTH⟩ ⟨FREQUENCY⟩ | 2500 : 55 to 3000 : 75 | A |
| | AddDynNoise | It adds a dynamic white noise to the watermarked signal | ⟨STRENGTH⟩ | 20 to 40 | B |
| | AddFFTNoise | It adds white noise to the watermarked signal in the frequency domain | ⟨FFTSIZE⟩ ⟨STRENGTH⟩ | 256 : 1000 to 1024 : 3000 | C |
| | AddNoise | A white Gaussian noise is contaminated the watermarked signal to simulate ambient distortion | ⟨STRENGTH⟩ | 35 dB level to 5 dB | D |
| | AddSinus | It adds a sinus signal to the watermarked signal | ⟨AMPLITUDE⟩ ⟨FREQUENCY⟩ | 120 : 3000 to 150 : 3500 | E |
| Conversion | Resampling | The sampling rate of the watermarked signal is converted to ⟨SAMPLERATE1⟩ and then is reconverted to ⟨SAMPLERATE2⟩ | ⟨SAMPLERATE1⟩ ⟨SAMPLERATE2⟩ | 4 KHz : 16 KHz to 8 KHz : 16 KHz | F |
| | Requantization | The sample of the watermarked signal is quantized to ⟨QUANTIZATION1⟩ and then is requantized to ⟨QUANTIZATION2⟩ | ⟨QUANTIZATION1⟩ ⟨QUANTIZATION2⟩ | 8 bits and 16 bits | G |
| | Invert | It inverts all samples in the watermarked signal, like a 180 degree phase shift | NO PARAMETER REQUIRED | None | H |
| Ambience | Echo | An echo with a delay ⟨DELAY⟩ and decay ⟨DECAY⟩ is added to the watermarked signal | ⟨DELAY⟩ ⟨DECAY⟩ | 20 ms and 10% to 100 ms and 50% | I |
| Sample permutations | Cut samples | ⟨REMOVENUMBER⟩ samples are removed from the watermarked signal from every ⟨REMOVEDIST⟩ period | ⟨REMOVEDIST⟩ ⟨REMOVENUMBER⟩ | 1 and 1000 to 7 and 1000 | J |
| | Copy samples | Some of the samples of the watermarked signal are copied between the samples values | ⟨PERIOD⟩ ⟨COPYDIST⟩ ⟨COPYCOUNT⟩ | 1000 : 100 : 30 to 1000 : 200 : 60 | K |
| | LSB Zero | Set all samples of the watermarked signal to zero | NO PARAMETER REQUIRED | None | L |
| | Smooth | The new sample value depends on the samples before and after the modifying point | NO PARAMETER REQUIRED | None | M |
| | Stat1 | It averages the sample with its next neighbors | NO PARAMETER REQUIRED | None | N |
| Dynamics | Amplify | The amplitude of the watermarked signal is increased up to ⟨FACTOR1⟩ and is decreased down to ⟨FACTOR2⟩, respectively | ⟨FACTOR1⟩ ⟨FACTOR2⟩ | 150% and 75% 200% and 50% | O |
| | Denoising | The watermarked signal is denoised by ⟨FACTOR⟩ | ⟨FACTOR⟩. | −80 dB to −60 dB | P |
| Filters | Low Pass Filter (LPF) | The watermarked signal is filtered by an elliptic LPF with cutoff frequency of ⟨FREQUENCY⟩ | ⟨FREQUENCY⟩ | 5 KHz to 4 KHz | Q |
| | Band Pass Filter (BPF) | The watermarked signal is filtered by an elliptic filter with bandwidth from ⟨FREQUENCY1⟩ to ⟨FREQUENCY2⟩ to simulate a narrowband telephony channel | ⟨FREQUENCY1⟩ ⟨FREQUENCY2⟩ | 500 Hz & 4000 Hz to 300 Hz & 3400 Hz | R |
| | High Pass Filter (HPF) | The watermarked signal is filtered by an elliptic HPF with cutoff frequency of ⟨FREQUENCY⟩ | ⟨FREQUENCY⟩ | 500 Hz to 800 Hz | S |

Table 1: Continued.

| Attack type | Attack name | Description | Parameter(s) | Default value(s) | |
|---|---|---|---|---|---|
| Time stretch and pitch shift | Pitch scale | The pitch of the watermarked signal is nonlinearly scaled without changing the time | ⟨SCALEFACTOR⟩ | 1.05 to 1 : 10 | T |
| | Time stretch | The time of the watermarked signal is nonlinearly stretched | ⟨TEMPOFACTOR⟩ | 1.05 to 1.10 | U |
| Compression | CELP coding | The watermarked signal is codded with rate of ⟨BITRATE⟩ by CELP codecs and then is decoded to original one | ⟨BITRATE⟩ | 16 Kbps to 9.6 kbps | V |
| | MP3 compression | The watermarked signal is compressed by MP3 with different rate ⟨BITRATE⟩ | ⟨BITRATE⟩ | 128 to 32 | W |
| | G.711 | The watermarked signal is codded by standard 64 kbps, A/$\mu$-law PCM | NO PARAMETER REQUIRED | None | X |



Figure 8: Variation of the BSC capacity with respect to different WNRs for different quantization steps.
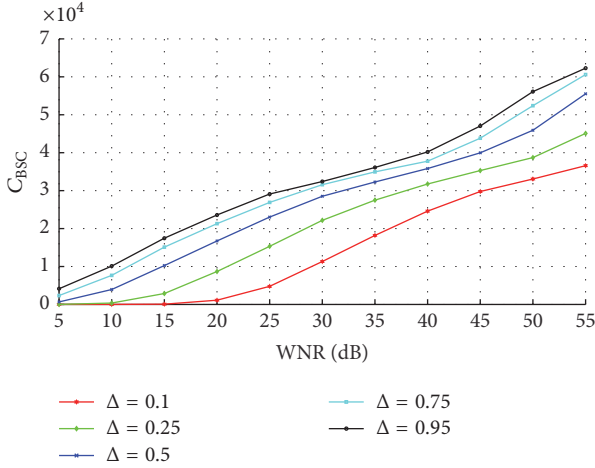


Figure 9: Variation of the BSC capacity with respect to different WNRs for different frame lengths.

Table 2 compares the BER with state-of-the-art speech watermarking techniques. We implemented all the techniques and tested them for the entire TIMIT corpus under different attacks. As can be observed, the proposed speech watermarking technique has a lower BER overall compared with other techniques.

The perceptual quality of the watermarked signal is critical for the evaluation of the proposed watermarked technique, which can be measured based on the mean opinion score (MOS) (as proposed by the International Telecommunications Union (ITU-T) [23]) and SNR. The MOS uses a subjective evaluation technique to score the watermarked signal, which is presented in Table 3. In the MOS evaluation method, 10 people were asked to listen blindly to the original and watermarked signals. Then they reported the dissimilarities between the quality of the original and watermarked speech signals. The average of these reports were computed for MOS music and MOS speech and presented in Table 4.
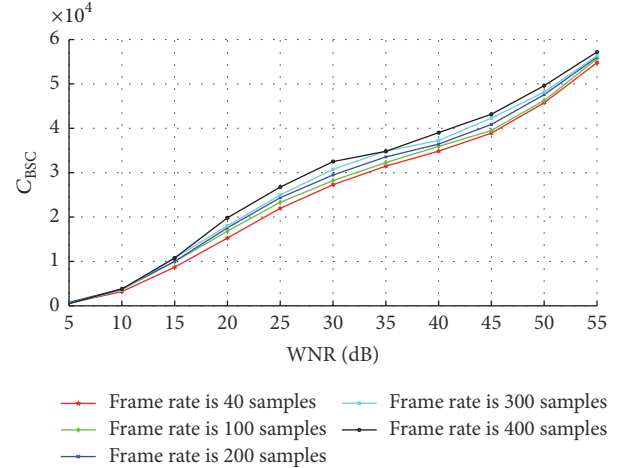
An objective evaluation technique, such as SWR and SNR, attempts to quantify this amount based on the following formula:

$$\text{SNR} = 10 \times \log_{10} \frac{\sum_n S^2}{\sum_n \left(\widehat{S} - S\right)^2}, \tag{31}$$

where $S$ and $\widehat{S}$ are the original and watermarked signals, respectively.

Table 4 presents a comparison of the proposed technique and other techniques in terms of imperceptibility and capacity. Based on the results, it seems that the proposed speech watermarking technique outperformed the other techniques in terms of capacity and imperceptibility. Although the SNR for formant tuning [21] is higher than the proposed technique, the capacity and robustness of the proposed technique are greater than those for formant tuning [21] and Analysis-by-Synthesis [22].

TABLE 2: Comparison with the robustness of different speech watermarking techniques in terms of BER (%).

| Attack | The proposed method | DWPT+ multiplication [14] | Formant tuning [21] | Analysis-by-Synthesis [22] |
|---|---|---|---|---|
| No attack | 0.00 | 0.00 | 0.04 | 0.06 |
| A | 1.91–4.23 | 2.09–5.43 | 3.65–6.45 | 7.96–9.65 |
| B | 9.65–21.77 | 10.45–22.32 | 12.76–24.45 | 16.23–25.23 |
| C | 10.13–20.43 | 12.43–21.32 | 14.23–23.54 | 17.43–26.32 |
| D | 10.53–19.23 | 10.33–18.93 | 11.63–23.23 | 15.33–25.98 |
| E | 0.32–2.02 | 0.763–1.14 | 1.23–2.32 | 2.98–4.32 |
| F | 13.54–17.23 | 14.32–17.65 | 26.23–37.83 | 29.45–33.06 |
| G | 3.23 | 2.65 | 19.32 | 23.87 |
| H | 0.23 | 0.00 | 9.43 | 12.45 |
| I | 1.34–4.65 | 2.34–5.11 | 4.65–10.43 | 8.23–16.43 |
| J | 1.23–2.54 | 1.32–4.67 | 6.54–10.54 | 11.54–18.87 |
| K | 1.32–3.16 | 1.78–4.23 | 7.51–10.34 | 11.49–19.43 |
| L | 0.92 | 1.98 | 1.50 | 4.04 |
| M | 3.12 | 5.76 | 10.34 | 21.68 |
| N | 4.10 | 4.23 | 6.65 | 9.54 |
| O | 1.21–2.54 | 0.00–1.43 | 5.97–8.76 | 8.98–15.54 |
| P | 1.00–3.54 | 2.43–5.43 | 9.65–14.56 | 19.65–26.45 |
| Q | 21.43–29.43 | 24.54–31.43 | 40.54–44.43 | 50.09–50.32 |
| R | 4.84–9.54 | 5.32–10.32 | 16.65–29.44 | 20.54–36.98 |
| S | 13.32–18.54 | 15.00–19.43 | 20.43–29.23 | 28.54–30.76 |
| T | 1.32–2.32 | 2.01–3.13 | 7.43–10.43 | 9.65–15.32 |
| U | 0.15–0.23 | 0.18–0.43 | 1.45–3.21 | 4.32–5.43 |
| V | 6.54–9.54 | 11.43–14.54 | 1.32–4.21 | 2.32–4.32 |
| W | 10.43–20.34 | 11.43–25.34 | 36.32–45.65 | 33.43–50.32 |
| X | 23.11 | 24.17 | 48.32 | 50.65 |
| Average | 5.80–9.04 | 6.68–10.04 | 12.95–17.39 | 16.82–21.48 |

TABLE 3: MOS grades [23].

| MOS | Quality | Quality scale | Effort required to understand meaning scale |
|---|---|---|---|
| (5) | Excellent | Imperceptible | No effort required |
| (4) | Good | Perceptible, but not annoying | No appreciable effort required |
| (3) | Fair | Slightly annoying | Moderate effort required |
| (2) | Poor | Annoying | Considerable effort required |
| (1) | Bad | Very annoying | No meaning was understood |

As observed in Table 4, each entity was bounded between two values that related a particular value of imperceptibility (SNR and MOS) to a particular capacity. Consequently, when the capacity increased, imperceptibility decreased. The trade-off value is completely application dependent and should be determined by the user.

## 5. Performance Analysis

Generally, two types of errors, false positive probability (FPP) and false negative probability (FNP), must always be analyzed to validate the security of a watermarking system [25]. FPP is defined when an unwatermarked speech signal is declared as a watermarked speech signal by the watermark extractor. Similarly, FNP is defined when the watermarked speech signal is declared as an unwatermarked speech signal by the

watermark extractor. By assuming that the watermark bits are independent random variables, both the FPP and FNP can be formulated based on Bernoulli trials, which is expressed as follows:

$$
P_e = \underbrace{\sum_{i=0}^{T-1} \binom{N}{i} P_{\text{FN}}^i \left(1 - P_{\text{FN}}\right)^{(N-i)}}_{\text{FNP}}
$$
$$
+ \underbrace{\sum_{i=T}^{N} \binom{N}{i} P_{\text{FP}}^i \left(1 - P_{\text{FP}}\right)^{(N-i)}}_{\text{FPP}},
\tag{32}
$$

where $N$ is the total number of watermark bits; $i$ is the number of matching bits; $\binom{N}{i}$ is a binomial coefficient; $P_{\text{FP}}$ is the probability of a false positive, which is assumed to be 0.5;

TABLE 4: Comparison of various watermarking techniques in terms of payload and imperceptibility.

| Technique | Quality scale | Effort required to understand meaning scale | SNR (dB) | Theoretical payload (bps) |
|---|---|---|---|---|
| Analysis-by-Synthesis [22] | 4.01–3.80 | 4.76–3.95 | 28.08–25.32 | 33.33–50 |
| Formant tuning [21] | 4.98–4.32 | 5.00–4.55 | 30.32–27.54 | 33.33–50 |
| DWPT+ multiplication [14] | 4.32–3.10 | 5.00–3.55 | 37.21–20.08 | 31.25–125 |
| The proposed method | 4.87–3.65 | 5.00–4.05 | 42.11–20.71 | 40–400 |



FIGURE 10: FPP with respect to various total number of watermark bits for different BER.
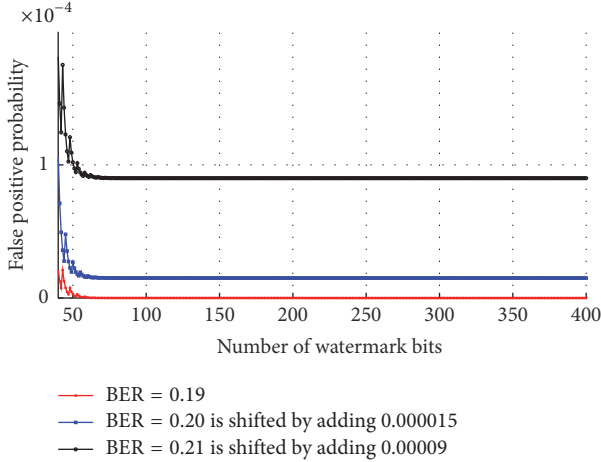


FIGURE 11: FNP with respect to various total number of watermark bits for different BER.

$P_{FN}$ is the probability of a false negative, which is assumed to be 0.0919 (as in Table 2); and $T$ is the threshold, which is computed as follows:

$$T = \lceil (1 - BER) \times N \rceil. \tag{33}$$

Figure 10 shows the FPP with respect to various total number of watermark bits for different BER. For better visualization, each line was shifted by adding a constant. As observed, the FPP was close to zero for $N$ greater than 50. There was a small fluctuation for $N$ less than 50, which depended on the BER.

Figure 11 shows the FNP with respect to various total number of watermark bits for different BER. For better visualization, each line was shifted by adding a constant. As can be observed, the FNP was close to zero for $N$ greater than 100. Additionally, whenever the BER decreased, the fluctuation increased.

## 6. Conclusion and Future Work

In this paper, a gain invariant speech watermarking technique was developed using the Lagrange optimization method. For this purpose, samples of the signal were separated based on odd and even indices. Then the ratio between the Lp-norms was quantized using the QIM method. Finally, the Lagrange method was used to estimate the optimized values. In a similar manner, the extraction process detected the watermark data blindly by finding the nearest quantization step.
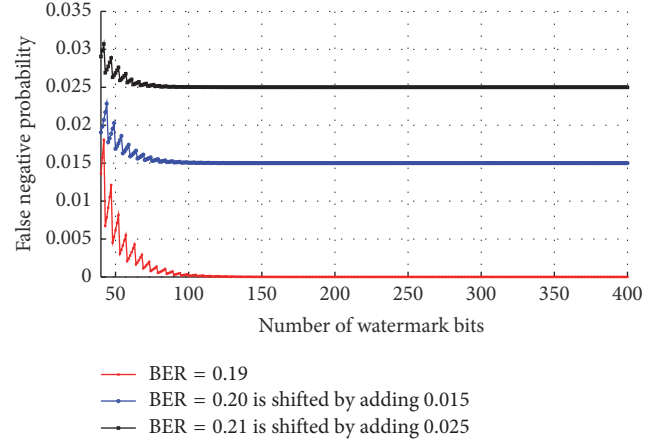
By assuming Laplacian distribution for the speech signal and Gaussian distribution for the noise signal, the probability of error and watermarking distortion were modeled based on a statistical analysis of the proposed technique. Additionally, experimental results not only proved that the developed watermarking technique was highly robust against different attacks, such compression, AWGN, filtering, and resampling, but also demonstrated the validity of the analytical model. For future work, an investigation on synchronization and adaptive quantization techniques might contribute to the proposed watermarking technique.

## Appendix

## A. Estimation of the Mean and Variance of the Ratio of Two Laplacian Variables Based on Taylor Series

In [26], the bivariate second-order Taylor expansion for $f(x, y)$ around $\theta = (E(x), E(y))$ is expressed as follows:

$$
\begin{aligned}
f(x, y) = & f(\theta) + f'_x(\theta)(x - \theta_x) + f'_y(\theta)(y - \theta_y) \\
& + \frac{1}{2}\left\{ f''_{xx}(\theta)(x - \theta_x)^2 \right. \\
& \left. + 2 f''_{xy}(\theta)(x - \theta_x)(y - \theta_y) + f''_{yy}(\theta)(y - \theta_x)^2 \right\} \\
& + \text{remainder}.
\end{aligned} \tag{A.1}
$$

Therefore, $E[f(X, Y)]$ can be expanded about $\theta = (E(X), E(Y))$ to compute the approximate values as follows:

$$
\begin{aligned}
E(f(X, Y)) = f(\theta) + \frac{1}{2} &\Big\{ f''_{xx}(\theta) \operatorname{var}(X) \\
&+ 2 f''_{xy}(\theta) \operatorname{cov}(X, Y) + f''_{yy}(\theta) \operatorname{var}(Y) \Big\} \\
&+ O(n^{-1}).
\end{aligned} \tag{A.2}
$$

For $f = R/S$, $f''_{RR} = 0$, $f''_{RS} = -S^{-2}$, and $f''_{SS} = 2R/S^3$. Then, the mean and variance of the ratio between $R$ and $S(E(R/S))$, respectively, can be estimated as follows:

$$
\begin{aligned}
E\left(\frac{R}{S}\right) &\equiv E(f(R, S)) \\
&\approx \frac{E(R)}{E(S)} - \frac{\operatorname{cov}(R, S)}{E(S)^2} + \frac{\operatorname{var}(S) E(R)}{E(S)^3} \\
&= \frac{\mu_R}{\mu_S}\left(1 + \frac{\sigma_S^2}{\mu_S^2}\right),
\end{aligned}
$$

$$
\begin{aligned}
\operatorname{var}\left(\frac{R}{S}\right) &\approx \frac{1}{E^2 S} \operatorname{var}(R) + 2 \frac{-ER}{E^3 S} \operatorname{cov}(R, S) \\
&\quad + \frac{E^2 R}{E^4 S} \operatorname{var}(S) \\
&= \frac{\mu_R^2}{\mu_S^2}\left[\frac{\sigma_R^2}{\mu_R^2} - 2 \frac{\operatorname{cov}(R, S)}{\mu_R \mu_S} + \frac{\sigma_S^2}{\mu_S^2}\right] \\
&= \frac{\mu_R^2}{\mu_S^2}\left(\frac{\sigma_S^2}{\mu_S^2} - \frac{\sigma_S^4}{\mu_S^4}\right) + \frac{\sigma_R^2}{\mu_S^2}.
\end{aligned} \tag{A.3}
$$

## B. Compute the Absolute Moment of the Laplacian Distribution

The moment of Laplacian distribution expressed as follows:

$$
\begin{aligned}
E(|X|^n) &= \int_{-\infty}^{\infty} |X|^n \cdot \frac{1}{2b} \cdot e^{-((X-\mu)/b)} dx \\
&= \frac{1}{2b} \int_{-\infty}^{\infty} |X|^n \cdot e^{-((X-\mu)/b)} dx.
\end{aligned} \tag{B.1}
$$

There are two cases, $X \geq \mu$ and $X < \mu$:

$$
E(|X|^n)
= \begin{cases}
\text{If } X \geq \mu \quad \text{then } \dfrac{1}{2b} \displaystyle\int_{-\infty}^{\infty} |X|^n \cdot e^{-((X-\mu)/b)} dx \\[2mm]
\text{If } X < \mu \quad \text{then } \dfrac{1}{2b} \displaystyle\int_{-\infty}^{\infty} |X|^n \cdot e^{-((\mu-X)/b)} dx.
\end{cases} \tag{B.2}
$$

For first case, when $X \geq \mu$,

$$
\begin{aligned}
E(|X|^n) = \frac{1}{2b} \Bigg[ &\int_{-\infty}^{0} -X^n \cdot e^{-((X-\mu)/b)} dx \\
&+ \int_{0}^{\infty} X^n \cdot e^{-((X-\mu)/b)} dx \Bigg]
\end{aligned}
$$

$$
= \frac{e^{\mu/b}}{2b} \left[ (-1)^n \underbrace{\int_{-\infty}^{0} X^n e^{X/b} dx}_{I_n} + \underbrace{\int_{0}^{\infty} X^n e^{-X/b} dx}_{I} \right]. \tag{B.3}
$$

If $t = -X/b$, then $I$ can be expressed as

$$
I = b^{n+1} \int_{0}^{\infty} t^n e^{-t} dt = b^{n+1} \cdot n! = n!, \tag{B.4}
$$

$I_n$ can also be expressed as

$$
\begin{aligned}
I_n &= \int_{-\infty}^{0} X^n e^{X/b} dx = \int_{-\infty}^{0} (b \cdot t)^n e^{-t} \cdot b \cdot dt \\
&= b^{n+1} \int_{-\infty}^{0} t^n e^{-t} dt \\
&= \left. \frac{t^n e^{-t}}{-1} \right|_{-\infty}^{0} - \int_{-\infty}^{0} \frac{n \cdot t^{n-1} e^{-t}}{-1} dt = 0 + n I_{n-1}.
\end{aligned} \tag{B.5}
$$

Substituting (B.4) and (B.5) into (B.3), the absolute moment of the Laplacian distribution can be computed based on

$$
E(|X|^n) = \left(\frac{e^{\mu/b} b^n}{2}\right) \left[(-1)^n \cdot I_n + n!\right]. \tag{B.6}
$$

## Competing Interests

The authors declare that they have no competing interests.

## References

[1] M. A. Nematollahi, C. Vorakulpipat, and H. G. Rosales, *Digital Watermarking: Techniques and Trends*, vol. 11, Springer, 2016.

[2] M. A. Nematollahi and S. A. R. Al-Haddad, "An overview of digital speech watermarking," *International Journal of Speech Technology*, vol. 16, no. 4, pp. 471–488, 2013.

[3] H.-T. Hu and L.-Y. Hsu, "A DWT-based rational dither modulation scheme for effective blind audio watermarking," *Circuits, Systems, and Signal Processing*, vol. 35, no. 2, pp. 553–572, 2016.

[4] B. Chen and G. W. Wornell, "Quantization index modulation: a class of provably good methods for digital watermarking and information embedding," *Institute of Electrical and Electronics Engineers. Transactions on Information Theory*, vol. 47, no. 4, pp. 1423–1443, 2001.

[5] M. A. Nematollahi, M. A. Akhaee, S. A. R. Al-Haddad, and H. Gamboa-Rosales, "Semi-fragile digital speech watermarking for online speaker recognition," *Eurasip Journal on Audio, Speech, and Music Processing*, vol. 2015, no. 1, article no. 31, 2015.

[6] P. Guccione and M. Scagliola, "Hyperbolic RDM for nonlinear valumetric distortions," *IEEE Transactions on Information Forensics and Security*, vol. 4, no. 1, pp. 25–35, 2009.

[7] N. Cai, N. Zhu, S. Weng, and B. Wing-Kuen Ling, "Difference angle quantization index modulation scheme for image watermarking," *Signal Processing: Image Communication*, vol. 34, pp. 52–60, 2015.

[8] X. Zhu and S. Peng, "A novel quantization watermarking scheme by modulating the normalized correlation," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP '12)*, pp. 1765–1768, IEEE, Kyoto, Japan, March 2012.

[9] M. A. Akhaee, S. M. E. Sahraeian, and C. Jin, "Blind image watermarking using a sample projection approach," *IEEE Transactions on Information Forensics and Security*, vol. 6, no. 3, pp. 883–893, 2011.

[10] N. K. Kalantari and S. M. Ahadi, "A logarithmic quantization index modulation for perceptually better data hiding," *IEEE Transactions on Image Processing*, vol. 19, no. 6, pp. 1504–1517, 2010.

[11] M. Zareian and H. R. Tohidypour, "A novel gain invariant quantization-based watermarking approach," *IEEE Transactions on Information Forensics and Security*, vol. 9, no. 11, pp. 1804–1813, 2014.

[12] M. Faundez-Zanuy, M. Hagmüller, and G. Kubin, "Speaker verification security improvement by means of speech watermarking," *Speech Communication*, vol. 48, no. 12, pp. 1608–1619, 2006.

[13] M. Faundez-Zanuy, M. Hagmüller, and G. Kubin, "Speaker identification security improvement by means of speech watermarking," *Pattern Recognition*, vol. 40, no. 11, pp. 3027–3034, 2007.

[14] M. A. Nematollahi, H. Gamboa-Rosales, M. A. Akhaee, and S. A. R. Al-Haddad, "Robust digital speech watermarking for online speaker recognition," *Mathematical Problems in Engineering*, vol. 2015, Article ID 372398, 12 pages, 2015.

[15] M. A. Nematollahi, H. Gamboa-Rosales, F. J. Martinez-Ruiz, J. I. de la Rosa-Vargas, S. A. R. Al-Haddad, and M. Esmaeilpour, "Multi-factor authentication model based on multipurpose speech watermarking and online speaker recognition," *Multimedia Tools and Applications*, pp. 1–31, 2016.

[16] M. A. Nematollahi, S. A. R. Al-Haddad, S. Doraisamy, and H. Gamboa-Rosales, "Speaker frame selection for digital speech watermarking," *National Academy Science Letters*, vol. 39, no. 3, pp. 197–201, 2016.

[17] S. Gazor and W. Zhang, "Speech probability distribution," *IEEE Signal Processing Letters*, vol. 10, no. 7, pp. 204–207, 2003.

[18] M. A. Akhaee, N. Khademi Kalantari, and F. Marvasti, "Robust audio and speech watermarking using Gaussian and Laplacian modeling," *Signal Processing*, vol. 90, no. 8, pp. 2487–2497, 2010.

[19] J. S. Garofolo and L. D. Consortium, *TIMIT: Acoustic-Phonetic Continuous Speech Corpus*, Linguistic Data Consortium, 1993.

[20] S. Verdú and T. S. Han, "A general formula for channel capacity," *IEEE Transactions on Information Theory*, vol. 40, no. 4, pp. 1147–1157, 1994.

[21] S. Wang and M. Unoki, "Speech watermarking method based on formant tuning," *IEICE Transactions on Information and Systems*, vol. 98, no. 1, pp. 29–37, 2015.

[22] B. Yan and Y.-J. Guo, "Speech authentication by semi-fragile speech watermarking utilizing analysis by synthesis and spectral distortion optimization," *Multimedia Tools and Applications*, vol. 67, no. 2, pp. 383–405, 2013.

[23] I. Rec, *P. 800: Methods for Subjective Determination of Transmission Quality*, International Telecommunication Union, Geneva, Switzerland, 1996.

[24] M. Steinebach, F. A. P. Petitcolas, F. Raynal et al., "StirMark benchmark: audio watermarking attacks," in *Proceedings of the International Conference on Information Technology: Coding and Computing*, IEEE, 2001.

[25] K. Vivekananda Bhat, I. Sengupta, and A. Das, "An audio watermarking scheme using singular value decomposition and dither-modulation quantization," *Multimedia Tools and Applications*, vol. 52, no. 2-3, pp. 369–383, 2011.

[26] R. C. Elandt-Johnson and N. L. Johnson, *Survival Models and Data Analysis*, Wiley Classics Library, John Wiley & Sons, New York, NY, USA, 1999.