# Farklı İçeriklere Sahip Wav Ses Dosyalarında Ki-kare Steganaliz Atağına Karşı Sağlamlık Değerlendirmesi

**Ali Durdu**[*1]

[*1] Ali DURDU Ankara Sosyal Bilimler Üniversitesi Siyasal Bilgiler Fakültesi Yönetim Bilişim Sistemleri, ANKARA

**Keywords**
Chi-square,
Data security,
Hidden data,
Steganalysis,
Steganography

**Abstract:** This paper presents the performance of the chi-square method for steganalysis purposes on audio files through various content. Predetermined rates of stego data have been embedded into 4500 distinct cover wav files via LSB (Least Significant Bit) method using three types of audio input files: music, human voice, animal sounds. Results of the experiment have revealed that hidden data in wav files containing dissimilar audio samples cannot be easily detected by chi-square steganalysis method. However, wav files containing similar audio samples demonstrate better performance compared to dissimilar samples.

## An Evaluation of the Robustness Chi-Square Steganalysis Method on Wav Audio Files with Various Content

**Anahtar Kelimeler**
Ki-kare,
Veri güvenliği,
Gizli veri,
Steganaliz,
Steganografi

**Öz:** Bu çalışmada, farklı içerikli wav ses dosyalarına LSB(Least Significant Bit) yöntemi ile gizlenen veriler ki-kare yöntemiyle steganaliz edilmiştir. Çalışmada müzik, insan ve hayvan seslerinden 3 kategoride 4500 farklı wav ses dosyasına farklı oranlarda veriler gizlenmiştir. Bu çalışmada ki-kare steganaliz yöntemi diğer çalışmalardan farklı olarak imgeler yerine wav ses dosyalarına uyarlanmıştır. Deney sonuçlarına göre benzersiz ses örnekleri (insan veya müzik ses örnekleri) içeren wav ses dosyalarına gizlenen verilerin ki-kare steganaliz yöntemi tarafından tespit edilemediği gözlenmiştir. Benzer ses örnekleri (hayvan ses örnekleri) içeren dosyalarda ise gizlenen veriler ki-kare steganaliz yöntemi tarafından tespit edilmiştir.

**\***İlgili Yazar, email: ali.durdu@asbu.edu.tr

## 1. Introduction

While development of the technology has provided fast communication applications, at the same time secure communication has come into prominence. Different methods have been developed for secure communication such as cryptography and steganography.

Steganography is a data hiding method in an innocent media to prevent initial observations from third parties. Thus, third parties are not aware of the presence of the secret data when steganography is concerned. Steganography can be applied to image, video, audio or text files. For example, Kuo et al. have proposed a novel high capacity data hiding algorithm which is based on multi-bit encoding function[1].

There are studies that are related to the finding of data which are hidden with the methods of steganography. These methods whose name is steganalysis tries to analyze the carrier file with the aid of different methods. In studies carried out, steganalysis methods in digital environment were generally applied on image files [2]–[15]. Wu et al. have proposed a frame of feature dimension reduction based semi-supervised learning for high-dimensional unbalanced JPEG image steganalysis [2]. Holub and Fridrich have introduced a novel feature set for steganalysis of JPEG images [4]. Nouri and Mansouri have proposed a new SVD-based feature set for steganalysis both in spatial and JPEG domains [5]. Mohammadi et al. have introduced a new feature-based blind steganalysis method [6]. In the literature, there are also steganalysis studies on audio files. [16]-[18]. Liu et al. have proposed

steganalysis method for detecting the presence of information-hiding behavior in wav audios. [16]. Ren et al. have proposed a method for detection of adaptive multirate (AMR) audio steganography [17]. Yavanoglu et al. have introduced an intelligent steganalysis method to investigate wave audio signals if they contain any steganographic content or not [18].

In this study, steganalysis has been carried out for audio files in content-based approach. In the study two data sets are used and in total 4500 audio files are steganalyzed. The hidden data have been embedded into the cover medium by LSB (Least Significant Bit) method. The imperceptibility variation of hidden data in the test files has been examined in respect to the content.

The steganography method which is used to conceal data to the audio files and chi-square steganalysis method which is used to identify the data are shown in the part 2. In part 3, data set which is used in the study and experimental results are presented. The result of the study is elucidated in part 4.

## 2. Material and Method

### 2.1. Steganography method

LSB substitution steganography method is used so that it can be an input for the steganalysis method in this study in order to conceal data to the data set which will be used. The reason for the preference of LSB steganography method is that chi-square steganalysis method is designed for the data which is concealed with LSB. In LSB method, each carrier is concealed into the last least significant bit of the every byte in the file. Hence, an 8-byte audio sample can accommodate a single secret data byte. The bits of secret data can be placed either successively or randomly in the bytes. The steganography algorithm prefers to locate the secret data in the last bits of each audio byte successively. The system has been evaluated with various amounts of inputs as the effect of size constraints can be eliminated by using randomly sized secret data.

The steganography algorithm initially reads the original audio file in wav format. The generated random secret message is embedded into an audio file at a desired data rate. The secret message is written in the last bits of the audio file and data has been hidden into the audio file at 100%. So as to hide data for the desired rate, the secret data is embedded then the rest is filled with the original audio file, and in this way stego file is created. The cover audio file is regenerated in wav format at the end of these procedures.

Audio samples can be 8, 16, 32 or 64 bits in length [19]. However, 16-bit wav audio files are widespread on the Internet. Steganography and steganalysis algorithms are compatible with all wav audio file formats regardless of the size. Hiding and detection procedures are carried out by steganography and steganalysis algorithms in accordance with the bit size of wav files that makes the system suitable for all wav formats. Incorrect changes in the data that are in the header section in Wav audio files might cause a complete disruption because the data in the header section represent the basic of the audio file. Therefore, a part of the header section cannot be used in data placement procedures.

### 2.2. Chi-square steganalysis method

Even though it is not possible to distinguish cover audio file from its original in terms of aural manner, it leaves some statistical marks on the file. Several steganalysis methods are available that analyze statistical marks which occurs after the concealment of data. For instance, global detection systems [26], [27] χ2 test (chi-square) [20], [21], RS analysis, histogram analysis [23], visual attacks [22], RQP methods [24] and JPEG steganalysis [25]. Some of these methods are used only for pictures; nevertheless, some of them can also be used for audio files. While in the studies of literature, chi-square steganalysis method is used on picture files, it is used on wav audio files in this study.

Significant distortion cannot be created by LSB method in a cover audio file and that distortions mostly cannot be detected by human ear. However, some statistical marks are left on the file. Chi-square method, which manipulates those marks, was discovered when Westfeld and Pfitzmann thought that LSBs in an image are not random. What they believed was that the occurrence of each of two pixels in each pair of values (PoV) tends to be far from the average of the PoV. In other words, the occurrence of even numbered pixels cannot possibly be close to the occurrence of odd numbered pixels in a stego image. It is observed that in Westfeld and Pfitzmann's studies the data in the original files, in this file the data is not concealed, odd and even numbered occurrences are not equal. Nevertheless, it is observed that occurrences in the file in which the data is concealed are equal [21]. Chi-square steganalysis method works by using this difference. When the data is concealed PoVs value in the original file alters. As a matter of fact, while embedding a secret data using LSB method, increment even numbered pixels by one or they remain unchanged. Odd values are increased once or there is no change [20]. As

a result of this change, each even in the histogram of pixel values approaches to the same number and the histogram of the file is "pair-wise".

In wav audio, concealed data at various rates are concealed with LSB method. Odd and even numbered occurrences are grouped in the same category. Differences between odd and even numbered occurrences are close to the zero when the file is full of concealed data. In each category, differences of odd and even numbered occurrences are calculated and the obtained histograms are demonstrated in the figure 1. Histograms show the data rate depending on the differences between the categories. To give an illustration, if concealed data carrier file is full of concealed data, difference between the categories is close to the zero.
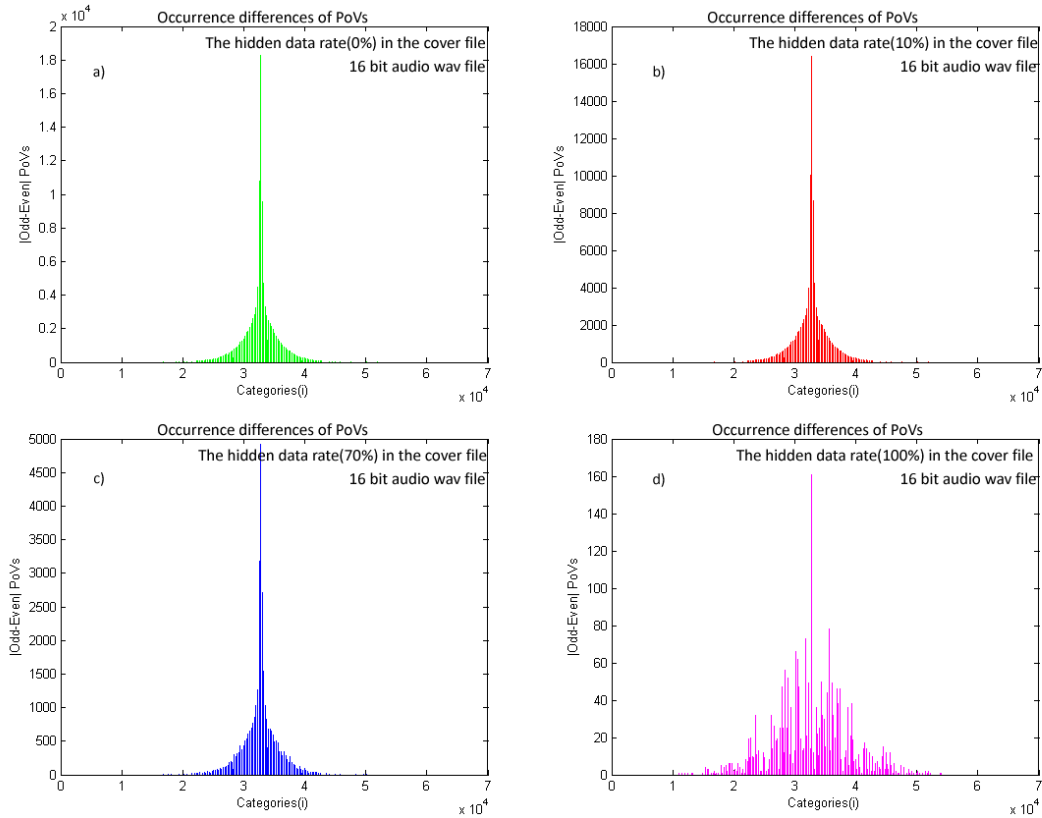


**Figure 1.** The differences of categorized odd and even occurrences of wav audio file with various rates of embedding message

## 3. Results

### 3.1. Audio sets

Audio Set 1: In this study, have used audio files incorporating music, human and animal voice files retrieved from YouTube [28]. These audio files have been split into 10-second samples using Audacity software [29] and resaved as 16-bit wav format. At the end of the process have obtained 3000 wav files, 1000 from each category.

Audio Set 2: For the second data set, 1500 files, 500 from each category i.e. music, human and animal voice have been downloaded from various servers [30, 31]. These audio files are, in 16-bit wav format, 10-seconds or shorter.

### 3.2. Content based chi-square analysis

Chi-square is an analysis method that works on odd and even numbered occurrence distributions. A sound sample comprising of odd and even numbered samples forms a PoVs. In this method, if the occurrence difference of PoVs is small i.e. close to zero, it is assumed that the content can incorporate hidden data. On the other hand (bigger than zero), the method returns a negative result. Thus, if an audio file contains dissimilar sounds, there will be too many sound samples in different occurrences. Therefore, each dissimilar sound sample constitutes a distinct category. Since there are so many different odd and even numbered occurrences, the number of categories will ascend and the difference between odd and even occurrences will decrease. This situation leads to false positive result in chi-square method. Some files can incorporate similar audio samples like animal

sounds. In this case, the number of odd and even numbered occurrences of audio samples is high. Hence, the categories pile on certain numbers and the number of distinct categories demonstrates lower values. In some cases, there are only a few categories but too many occurrences; in this case, the chi-square analysis generates true positive results for files with such content.

Figure 2 shows the result of chi-square analysis of audio files selected by three categories: animal, human and music sounds. For the experimental scenario, both similar and dissimilar audio samples from each category are selected. The chi-square analysis of animal audio file with similar sound samples with 0%, 30%, 50%, 70% and 100% hidden data produces true positive results as seen in Figure 2a. However, it generates false positives when animal audio files contain dissimilar sound samples as shown in Figure 2-b. While the method on human sound file containing similar audio samples embedded with hidden data rate of 0%, 30%, 50%, 70% and 100 % hidden data produces true positive results for 0 %, the results are partially true positive for hiding rate up to almost 12%, yet they are completely false positive beyond 12% as seen in Figure 2.c. The motive behind this peculiarity is that it is hard to find human sound samples with similar content. Even a person reciting the same sentence may not generate similar sound samples. In addition, music files with similar content reflex similar behavior. Hence, the method generates false positives beyond 20% at 0% hidden data rate and beyond 10% for the rest of hidden data rates as shown in Figure 2.e. Consequently, all files containing dissimilar sound samples, the method yields false positive as seen in Figure 2.b, 2.d, and 2.f.

Figure 3 demonstrates odd and even numbered occurrence differences (PoVs) of audio samples chosen from three categories, namely animal, human and music. Audio files with similar and dissimilar sound contents from each category have been selected for analysis. By examining the occurrence difference of PoVs of the animal sound file with similar sound samples with a hidden data rate of 0%, 30%, 50%, 70% and 100% in Figure 3.a, the remarkable difference can be observed. However, occurrence difference of PoVs of the animal sound file with dissimilar sound samples is insignificant as seen in Figure 3.b. This occurrence is also repeated for the rest of the categories. Therefore, similar sound samples create fewer categories compared to dissimilar ones. Additionally, a high number of sound samples in all categories leads to high number of occurrence of PoV. This fact is also valid for Figure 3.a and Figure 3.c, which is a direct consequence of chi-square analysis. It also explains why chi-square generates true positives for files that have similar sound samples. On the contrary, more categories appear for files that dissimilar sound samples and the number of occurrence in categories are minimal. That is why occurrence differences are also low as seen in Figure 3.b, 3.d, and 3.f. Apparently, obtained results conflicts with the outcomes of chi-square analysis. When analyzed a file containing hidden data, the number of occurrence difference of odd and even values in the file approaches to zero. On the other hand, although there is no hidden data in files with dissimilar sound samples, the numbers of occurrence differences are close to zero, which explains why chi-square analysis produces false positives for those files.

Figure 4 shows the relationship between chi-square sums (given in Eq-2) and the number of categories of audio files chosen from three categories, namely animal, human, and music. When chi-square sums are higher than the value of categories, chi-square method yields accurate results. As seen on Figure 4.a, when no hidden data in the cover files which are composed of similar audio samples, the chi-square sum is higher than the value of categories. On the contrary, when the cover files include dissimilar audio samples with the same hidden data rate the value of categories are higher than chi-square sums as seen in Figure 4.b. This is because of the fact that fewer categories are formed in similar audio samples. However, the numbers of occurrences of audio samples in categories are increased. So, the more occurrences the categories have the more chi-square sums obtain. In Figure 4.a, 4.c and 4.e, chi-square method has produced accurate results, while in Figure 4.b, 4.d and 4.f it has produced inaccurate results. The reason for this is the low number occurrences in many categories and these causes to a low amount of chi-square sums. Since chi-square method utilizes the difference between odd and even values, it produces positive results because of a low level of frequency difference. However, such exceptional cases generate false positive results and therefore, hidden data on the files cannot be detected by chi-square method.

Audio Set 1: In this study has embedded hidden data into 3000 audio files with hidden data rates of 0%, 30%, 50%, 70% and 100 % followed by implementation of chi-square analysis on these files. As seen in Figure 5.a, 5.b, and 5.c chi-square method have produced false positive results when 3000 audio files with 10-second samples

are used. The main reason for having a massive rate of false positives is the low number of differences between odd and even samples as shown in Figure 6.a, 6.b, and 6.c. The number of categories in accordance with these results, outnumbered chi-square sums as seen in Figure 7.a, 7.b, and 7.c.
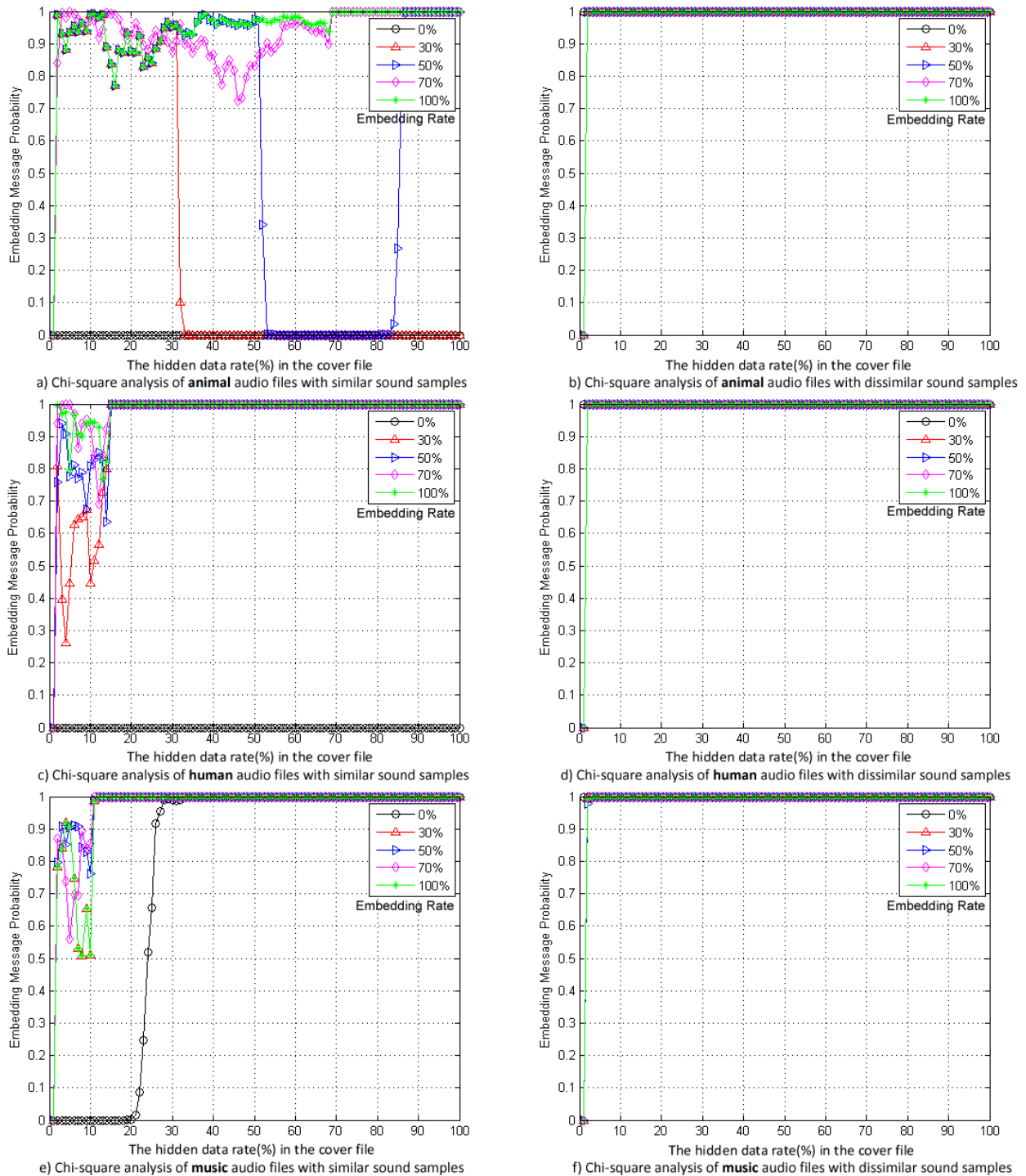


a) Chi-square analysis of **animal** audio files with similar sound samples

b) Chi-square analysis of **animal** audio files with dissimilar sound samples

c) Chi-square analysis of **human** audio files with similar sound samples

d) Chi-square analysis of **human** audio files with dissimilar sound samples

e) Chi-square analysis of **music** audio files with similar sound samples

f) Chi-square analysis of **music** audio files with dissimilar sound samples

**Figure 2.** Chi-square analysis of animal, human and music audio files containing similar and dissimilar sound sample

Audio Set 2: In another experiment, has employed 1500 audio files, which have a lower number of categories, to embed hidden data with rates of 0%, 30%, 50%, 70%, and 100%. After that chi-square analysis has been applied on these files. Figure 8, 9 and 10 shows the average values obtained from the analysis. Results seen in Figure 5, 6, 7 and Figure 8, 9, 10 support each other. However, the results obtained from the analysis for Audio Set 2 are more applicable than those for Audio Set 1. This is due to the fact that categories can occur and this is essentially compatible with the operating logic of chi-square method.
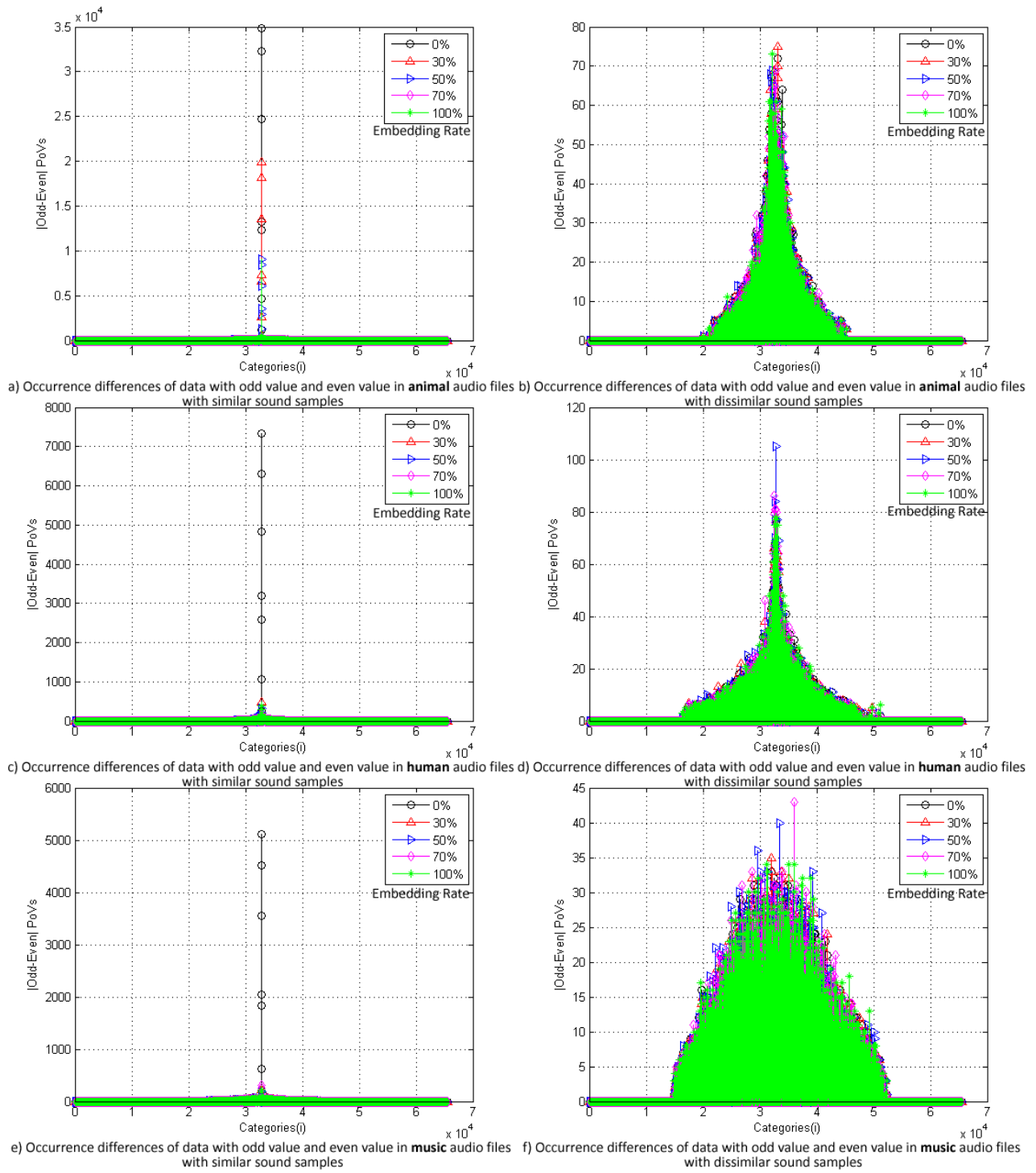
a) Occurrence differences of data with odd value and even value in **animal** audio files with similar sound samples

b) Occurrence differences of data with odd value and even value in **animal** audio files with dissimilar sound samples

c) Occurrence differences of data with odd value and even value in **human** audio files with similar sound samples

d) Occurrence differences of data with odd value and even value in **human** audio files with dissimilar sound samples

e) Occurrence differences of data with odd value and even value in **music** audio files with similar sound samples

f) Occurrence differences of data with odd value and even value in **music** audio files with dissimilar sound samples

**Figure 3.** Occurrence differences of odd and even numbered audio samples in animal, human and music audio files with similar and dissimilar audio samples with different rates of embedding messages
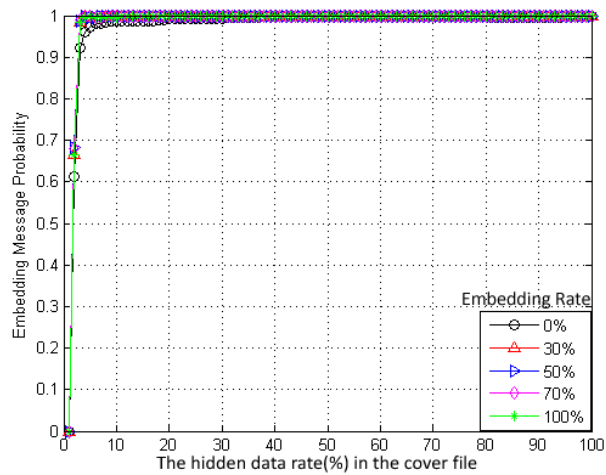
a) Categories and chi-square sums of **animal** audio files with similar sound samples  b) Categories and chi-square sums of **animal** audio files with dissimilar sound samples

c) Categories and chi-square sums of **human** audio files with similar sound samples  d) Categories and chi-square sums of **human** audio files with dissimilar sound samples

e) Categories and chi-square sums of **music** audio files with similar sound samples   f) Categories and chi-square sums of **music** audio files with dissimilar sound samples

**Figure 4.** Categories and chi-square sums of animal, human and music audio files with similar and dissimilar sound samples with various embedding rates

a) Average chi-square graph analyses of 1000 wav **animal** audio files
with different rates of embedding a messages



b) Average chi-square graph analyses of 1000 wav **human** audio files
with different rates of embedding a messages



c) Average chi-square graph analyses of 1000 wav **music** audio files
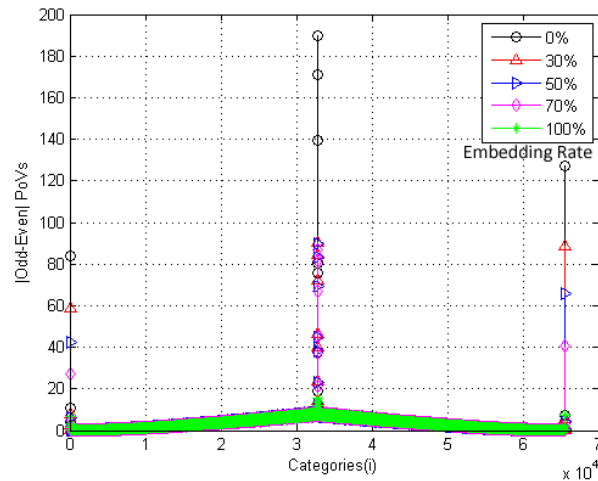with different rates of embedding a messages

**Figure 5.** Audio Set 1: Average chi-square graph analyses of 3000 wav different type (animal, human and music) audio files
Audio Set 2 has shorter audio files and thus, fewer

185

a) Average occurence differences graph of data with odd  and even value analyses
of 1000 wav **animal** audio files with different rates of embedding  a messages

b) Average occurence differences graph of data with odd  and even value analyses
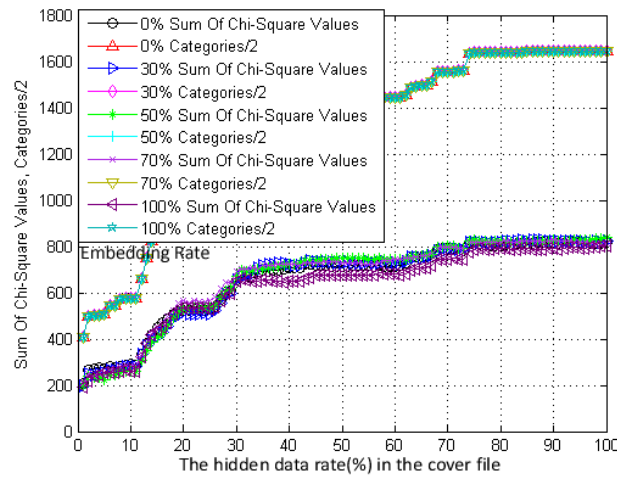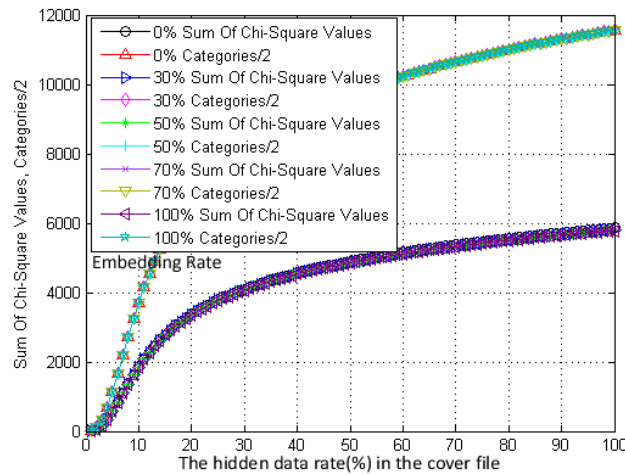of 1000 wav **human** audio files with different rates of embedding  a messages

c) Average occurence differences graph of data with odd  and even value analyses
of 1000 wav **music** audio files with different rates of embedding  a messages

**Figure 6.** Audio Set 1: Average occurrence differences graph of data analyses of 3000 wav different type (animal, human music) audio files

a) Average graph categories and chi-square sums analyses
of 1000 wav **animal** audio files with different rates of embedding a messages
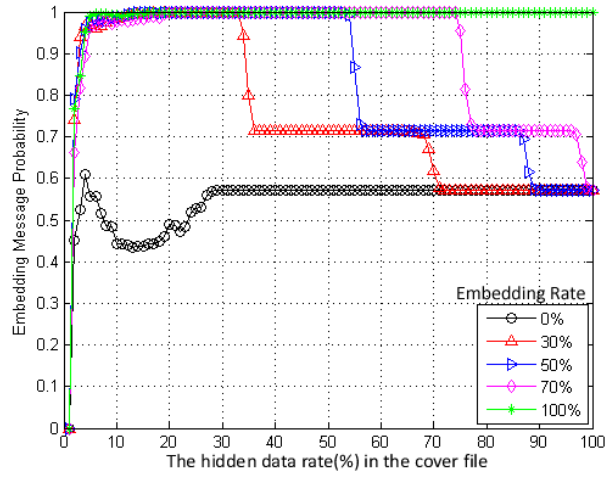


b) Average graph categories and chi-square sums analyses
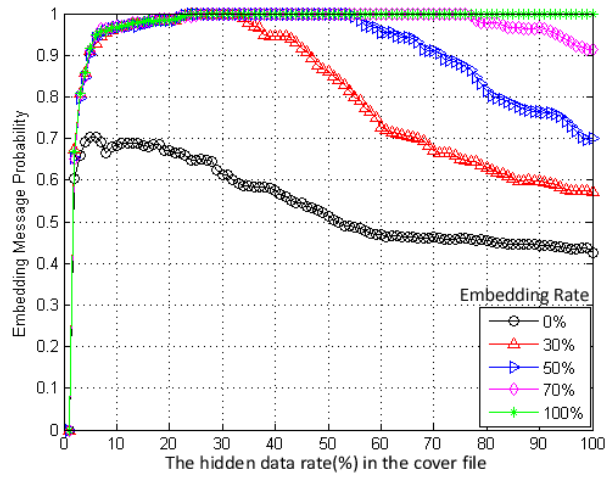of 1000 wav **human** audio files with different rates of embedding a messages



c) Average graph categories and chi-square sums analyses
of 1000 wav **music** audio files with different rates of embedding a messages
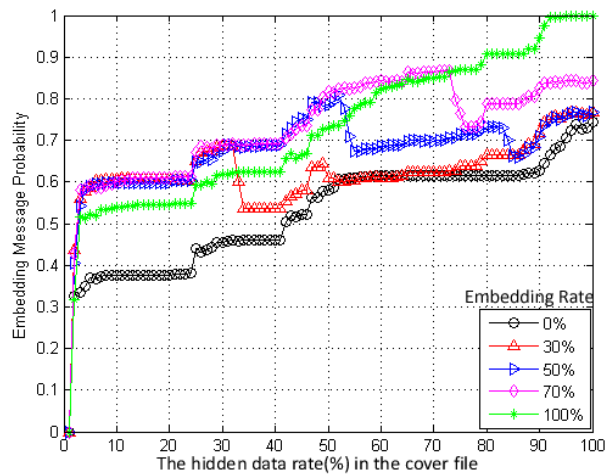
**Figure 7.** Audio Set 1: Average graph categories and chi-square sums analyses of 3000 wav different type (animal, human and music) audio files

a) Average graph Chi-square analyses of 500 wav **animal** audio files
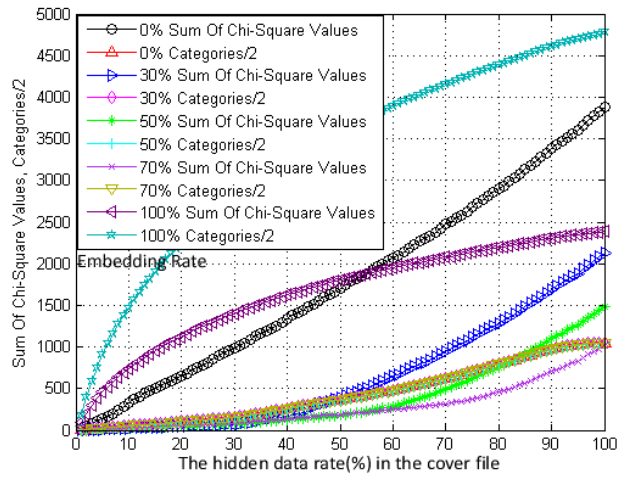with different rates of embedding a messages

b) Average graph Chi-square analyses of 500 wav **human** audio files
with different rates of embedding a messages
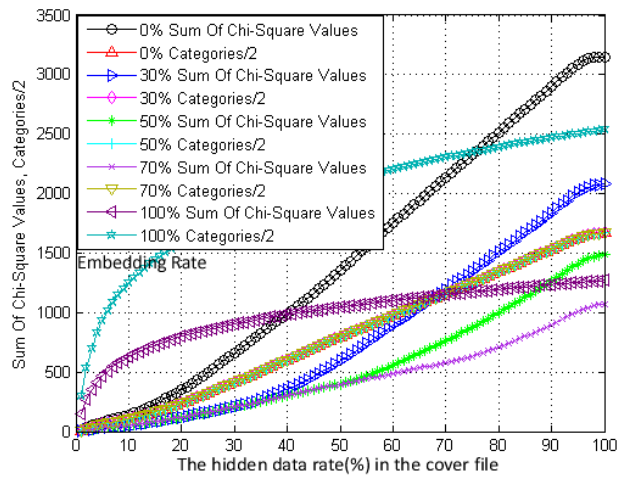
c) Average graph Chi-square analyses of 500 wav **music** audio files
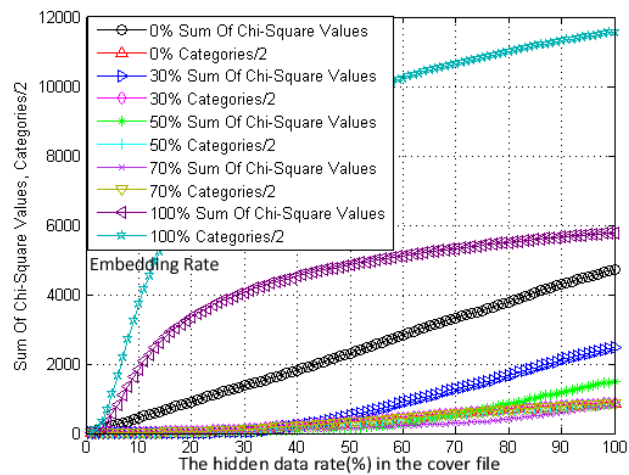with different rates of embedding a messages

**Figure 8.** Audio Set 2: Average chi-square graph analyses of 1500 wav different type (animal, human and music) audio files

188

a) Average graph categories and chi-square sums analyses
of 500 wav **animal** audio files with different rates of embedding a messages



b) Average graph categories and chi-square sums analyses
of 500 wav **human** audio files with different rates of embedding a messages



c) Average graph categories and chi-square sums analyses
of 500 wav **music** audio files with different rates of embedding a messages

**Figure 9.** Audio Set 2: Average occurrence differences graph of data analyses of 1500 wav different type (animal, human and music) audio files

a) Average graph occurrences differences of data with odd and even value analyses of 500 wav **animal** audio files with different rates of embedding a messages



b) Average graph occurrences differences of data with odd and even value analyses of 500 wav **human** audio files with different rates of embedding a messages



c) Average graph occurrences differences of data with odd and even value analyses of 500 wav **music** audio files with different rates of embedding a messages
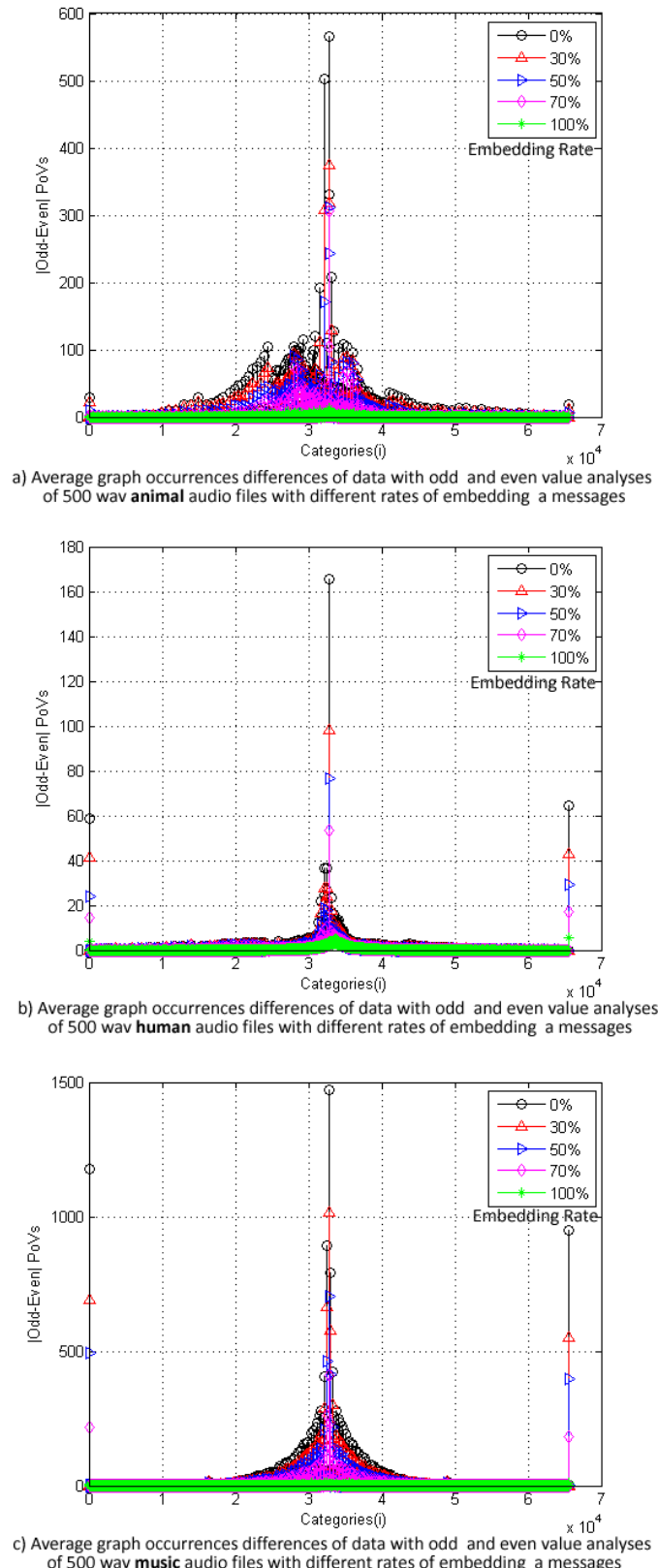
**Figure 10.** Audio Set 2: Average graph categories and chi-square sums analyses of 3000 wav different type (animal, human and music) audio files

When the embedding rate is 100% in musical audio files, it can produce 12000 different categories as seen in Figure 7.c and Figure 10.c. This implies more categories can be formed in musical audio files since they are rich in many dissimilar audio samples. While the difference between odd and even categories is 3000 in the case of 0% data hiding, it is only 550 in Figure 9.a. Comparing Figure 7.a, b, and c with Figure 10.a, b, and c respectively; it can be seen that the number of categories and chi-square sums remain the same at each hiding rate as audio samples with the same length. On the other hand, the number of categories and chi-square sums vary on Figure 8, 9 and 10 as it includes audio samples with various lengths. So we can conclude that the length of the sound file is of great influence on the performance of chi-square method. Hence, hiding data in shorter audio files results in better performance for chi-square method.

Some audio files contain similar sounds. Music files, for example, contain repetitive sounds in music tones of repeated parts and this increases the number of frequencies of odd and even audio samples, which then causes the categories to pile at certain numbers and also creates fewer different categories. In cases when there are only a few categories but too many frequencies, chi-square analysis produces true negative results. Thus, it is not safe to hide data in such files.

In conclusion, we can say that human sounds and music sounds have a high chance of containing unique sounds. It is extremely difficult to detect data hidden in these files with chi-square analysis. Yet, for audio files with similar sound samples, like animal sounds, chi-square produces more accurate results. It is understood from all experiments that data hiding in audio files with unique sound samples could not be detected by chi-square. In contrast, data hiding in audio files with similar sound samples could be partially detected by chi-square.

## 4. Discussion and Conclusion

In this paper, we have evaluated and analyzed the performance scores of four steganalysis methods on 16-bit audio wav files in which different rate of secret data embedded.

In our study, various rates of data have been hidden in 4500 different wav audio files from 3 categories of music, human and animal sounds and analyzed via chi-square analysis method. Analyses carried out reveal that chi-square produces proper results for audio files with similar sound samples while it produces inappropriate results for audio with dissimilar sound samples. Chi-square analysis produced accurate results for animal sounds with similar audio samples. It has, thus, become clear that preferring especially sound files, which include music and human voice for data hiding, will be safer so that they will not be detected by chi-square. Since the objective of the study is to evaluate and to analyze the chi-square method for steganography purposes using an extensive number of audio files, any comparison with the literature has not been provided.

## References

[1] Kuo, W.-C., Kuo, S.-H., Wang, C.-C., Wuu, L.-C. 2016. High capacity data hiding scheme based on multi-bit encoding function. Opt. - Int. J. Light Electron Opt, 127(2016), 4, 1762–1769.

[2] Wu, A., Feng, G., Zhang, X., Ren, Y. 2016. Unbalanced JPEG image steganalysis via multiview data match. J. Vis. Commun. Image Represent, 34(2016), 103–107.

[3] Li, M., Liu, Q. 2015. Steganalysis of SS Steganography: Hidden Data Identification and Extraction. Circuits, Syst. Signal Process, 34(2015), 10, 3305–3324.

[4] Holub, V., Fridrich, J. 2015, Low-complexity features for JPEG steganalysis using undecimated DCT. IEEE Trans. Inf. Forensics Secur., 10(2015), 2, 219–228.

[5] Nouri, R., Mansouri, A. 2015. Blind image steganalysis based on reciprocal singular value curve. in Machine Vision and Image Processing (MVIP), 2015 9th Iranian Conference on, 124–127.

[6] Mohammadi, F. G., Abadeh, M. S. 2014. Image steganalysis using a bee colony based feature selection algorithm. Eng. Appl. Artif. Intell., 31(2014), 35–43.

[7] Zhang, H., Ping, X. J., Xu, M. K., Wang, R. 2014, Steganalysis by subtractive pixel adjacency matrix and dimensionality reduction. Sci. China Inf. Sci., 57(2014), 4, 1–7.

[8] Lu, J. C., Liu, F. L., Luo, X. Y. 2014, Selection of image features for steganalysis based on the Fisher criterion, Digit. Investig., 11(2014), 1, 57–66.

[9] Pathak, P., Selvakumar, S. 2014, Blind Image Steganalysis of JPEG images using feature extraction through the process of dilation, Digit. Investig., 11(2014), 1, 67–77.

[10] Holub, V., Fridrich, J. 2013, Random projections of residuals for digital image steganalysis, IEEE Trans. Inf. Forensics Secur., 8(2013), 12, 1996–2006.

[11] Gul, G., Kurugollu, F. 2013, JPEG image steganalysis using multivariate PDF estimates with MRF cliques, IEEE Trans. Inf. Forensics Secur., 8(2013), 3, 578–587.

[12] Cho, S., Cha, B.-H., Gawecki, M., Jay Kuo, C.-C. 2013, Block-based image steganalysis: Algorithm and performance evaluation. J. Vis. Commun. Image Represent., 24(2013), 7, 846–856.

[13] Chen, G., Zhang, D., Zhu, W., Tao, Q., Zhang, C., Ruan, J. 2012, On optimal feature selection using harmony search for image steganalysis. in Proceedings - International Conference on Natural Computation, 1074–1078.

[14] Zong, H., Liu, F. L., Luo, X. Y. 2012, Blind image steganalysis based on wavelet coefficient correlation, Digit. Investig., 9(2012), 1, 58–68.

[15] Dai, Z., Xiong, Q., Peng, Y., Gao, H. 2012, Research on the large scale image steganalysis technology based on cloud computing and BP neutral network, in Proceedings of the 2012 8th International Conference on Intelligent Information Hiding and Multimedia Signal Processing, IIH-MSP 2012, 415–419.

[16] Liu Q., M., Sung, A. H., Qiao, M., 2008, Detecting information-hiding in WAV audios,. Pattern Recognition, 19th International Conference on, 1051-4651.

[17] Ren, Y. , Cai, T., Tang, M., Wang, L. 2015, AMR steganalysis based on the probability of same pulse position, IEEE Trans. Inf. Forensics Secur., 10(2015), 9, 1801–1811.

[18] Yavanoglu, U., Ozcakmak, B., Milletsever, O. 2012, A New Intelligent Steganalysis Method for Waveform Audio Files, 11th Int. Conf. Mach. Learn. Appl., Dec, 233–239.

[19] Microsoft Wave. 2018. http://soundfile.sapp.org/doc/WaveFormat (Accessed: 29.04.2018).

[20] Stanley, C. A. 2005, Pairs of Values and the Chi-squared Attack, 1–45.

[21] Westfeld, A., Pfitzmann, A. 2000, Attacks on Steganographic Systems Steganos, S-Tools and Some Lessons Learned, 1–16.

[22] Şahin, A. 2007, "New Methods on Image Steganography and Their Reliabilities, Trakya University, Natural and Applied Sciences, Ph.D. Thesis, Trakya.

[23] Yalman, Y. 2010, Design And Implementation Of A Steganography Method Based On Histogram Modification For Digital Images, Kocaeli University, Natural and Applied Sciences, Ph.D. Thesis, Kocaeli.

[24] Fridrich, J., Long, M. 2000, Steganalysis of LSB encoding in color images, IEEE Int. Conf. Multimed. Expo. ICME2000. Proceedings. Latest Adv. Fast Chang. World Multimed. 3(2000), c, 1279–1282.

[25] Fridrich, J., Goljan, M. 2002, Practical Steganalysis of Digital Images – State of the Art, Apr., 1–13.

[26] Farid, H. 2002, Detecting hidden messages using higher-order statistical models, in Proceedings. International Conference on Image Processing, 2, II-905-II-908.

[27] Farid, H. 2001, Detecting Steganographic Messages in Digital Images, Technical Report, Oct.

[28] Youtube, 2018. http://www.youtube.com (Accessed: 29.04.2018).

[29] Audacity, 2018. http://www.audacityteam.org (Accessed: 29.04.2018).

[30] Wav source, 2018. http://www.wavsource.com (Accessed: 29.04.2018).

[31] Daily Wav, 2018. http://www.dailywav.com (Accessed: 29.04.2018).