



**Scientific Computing Department
Faculty of Computer and Information Sciences
Ain Shams University**

Deoxyribonucleic Acid Based Steganography

Thesis submitted as a partial fulfillment of the requirements for the degree of
Master of Science in Computer and Information Sciences

By

Ghada Hamed Aly Kamel

Teaching Assistant at Scientific Computing Department,
Faculty of Computer and Information Sciences,
Ain Shams University

Under Supervision of

Prof. Dr. Mohamed Fahmy Tolba

Professor in Scientific Computing Department,
Faculty of Computer and Information Sciences,
Ain Shams University

Dr. Safaa Amin El-Sayed

Associate Professor in Scientific Computing Department,
Faculty of Computer and Information Sciences,
Ain Shams University

Dr. Mohammed Abd El-Rahman Marey

Assistant Professor in Scientific Computing Department,
Faculty of Computer and Information Sciences,
Ain Shams University

January– 2017
Cairo

Acknowledgment

First of all, I would like to thank GOD for his endless blessings, for giving me the power and strength to complete this work and for giving me supportive people.

Second, I would like to express my sincere gratitude to my supervisors; Prof. Dr. Mohamed Fahmy Tolba for his support, patience and guidance, Dr. Safaa Amin El-Sayed for her encouragement and Dr. Mohammed Marey for the special supervision experience he gave me. I am deeply thankful.

Third, I would like to thank my family for being available all the time and for the love they gave me through the years. Thank you for accepting me through the tough times and for always believing in me.

My dear friends who have helped me through the past time and kept on encouraging me to get this work done; Alaa Atef, Alaa Salah, Eman Hamdi, Hanan Yousry, Samar Ahmed and Yasmin Alaa without you it would have been much harder.

Last but not least, I would like to thank all my professors, colleagues and students who kept on encouraging me. Thank you for being in my life.

Ghada

Abstract

The research on developing data security techniques has been increased continuously, due to the indispensable need for powerful data protection in different applications, these include ownership protection, annotation, copyrighting, authentication and military applications. Cryptography and steganography are the most common and widely used techniques in computer security. Cryptography includes encrypting and decrypting some data into incomprehensible formats. While steganography aims to hide the existence of the message in a different media such as image, audio, video, etc. so that an intruder is mainly prevented from suspecting that the data is there. DNA is explored as a new carrier for data hiding due to its huge variations and its high capacity in which 10^6 TB of data can be stored in one gram of DNA. However, like every data storage device, hiding data in DNA requires a high protection level through a secure algorithm. This leads to a DNA steganography research field based on DNA computing, where the properties of biotechnology are exploited in steganography.

Mainly, there are two approaches in DNA based steganography topic. First approach: Techniques that provide high data hiding capacity but at the expense of the original biological properties of the used DNA reference sequence and it is the widest approach. This approach in most cases may result in dangerous side effects that may lead to the death of an organism related to the used sequence in the steganography process. The second approach cares about preserving the DNA biological and chemical properties

but at the expense of either the data hiding capacity allowed or the algorithms' cracking probability.

In this thesis, the second approach is achieved to allow merging the proposed algorithms with the industry in practical way and without any dangerous effects through two proposed algorithms. The first proposed one is a hybrid crypto-steganography algorithm is proposed that achieves double layers of security through two phases. In the first phase, the confidential data is converted to DNA using a proposed generic N-bits binary coding rule that leads to lower cracking probability compared by others. Then, DNA and amino acids Playfair is applied as a first phase to encrypt the DNA of the message resulting in ambiguity. In the second phase, the cipher text is placed with the ambiguity using 3:1 placement strategy. Then they are shuffled to be hidden in DNA at random positions generated by using a true real random number seed that is obtained from the atmospheric noise, thereby achieving very low cracking probability. The proposed technique is a blind preserved one as it achieves zero modification for the generated protein without extra data. The drawback of the first implemented algorithm in the steganography field is the low capacity which is overcome by merging the DNA based cryptography concepts with the proposed steganography technique that results in a double layer security level and low cracking probability. The achievement is that the selected cryptography algorithm results in the same size of the input data with no extra information.

The drawback of the first proposed algorithm gets a new objective which is providing higher data capacity using DNA with preserving its main structure. For the sake of this, the second proposed algorithm is implemented to combine the DNA conservative mutation concepts with the steganography. It results in utilizing the DNA in an effective and non-

negativity for neither the DNA nor the steganography process. The idea is that the conservative mutation is exploited in a way to support high capacity, thanks to the use of each DNA base to hide two bits that is considered the first achievement. Since, most of the steganography algorithms support high capacity at the expense of the DNA modification rate which prevents it from being carried out practically. The second achievement in this study is selecting the conservative mutation specially to solve the tradeoff between the high capacity feature and preserving the original biochemical properties of the used sequence which is one of the main challenges considered in this research work. Besides minimizing the data that's sent to the receiver to strengthen the algorithm's security as compared with other techniques, results in sending only the carrier DNA sequence containing the secret data which is the final and third achievement in the steganography field.

After some deep biological view, the second proposed algorithm is modified to generalize the conservative mutation for each amino acid by getting all its possible substitutions with the other amino acids that have the same structure and functions. This proposed modification results in minimizing the algorithm's cracking probability by nearly fourteen billion trials. Even in the worst-case scenario, if the intruder success in taking away some parts of data, the extracted part is not meaningful, since the data is hidden in an arbitrary DNA bases which complexes the cracking process.

Besides that, for the first-time real large sized data to three megabytes of different formats are used to evaluate the performance of the proposed algorithm and investigate its scalability. The experimental results prove that the proposed technique is the one that overcomes the weakness points in the current steganography techniques as it merges the highest capacity feature with preserving the biochemical properties of the used DNA sequence

through a blind and a highly-secure algorithm without any generated extra information and with the lowest achieved cracking probability.

List of Publications

- 1- Ghada Hamed, Mohammed Marey, Safaa Amin and Mohamed Fahmy Tolba. "DNA based steganography: Survey and analysis for parameters optimization." Applications of intelligent optimization in biology and medicine, Springer International Publishing, pp 47-89, 2016.
- 2- Ghada Hamed, Mohammed Marey, Safaa Amin and Mohamed Fahmy Tolba. "Hybrid technique for steganography-based on DNA with n-bits binary coding rule." Soft Computing and Pattern Recognition (SoCPaR), pp 95-102, IEEE, November 2015.
- 3- Ghada Hamed, Mohammed Marey, Safaa Amin and Mohamed Fahmy Tolba. "Hybrid Randomized and Biological Preserved DNA-Based Crypt-Steganography Using Generic N-Bits Binary Coding Rule". International Conference on Advanced Intelligent Systems and Informatics 2016, pp 618-627, Springer International Publishing, 2016.
- 4- Ghada Hamed, Mohammed Marey, Safaa Amin and Mohamed Fahmy Tolba. "Hybrid, Randomized and High Capacity Conservative Mutations DNA-Based Steganography for Large Sized Data". Information Sciences (Submitted).
- 5- Ghada Hamed, Mohammed Marey, Safaa Amin and Mohamed Fahmy Tolba. "Comparative Study for Various DNA Based Steganography Techniques with the Essential Conclusions about the Future

- Research”. The 11th IEEE International Conference on Computer Engineering and Systems, pp 220-225,IEEE, 2016.
- 6- Ghada Hamed, Mohammed Marey, Safaa Amin and Mohamed Fahmy Tolba. “Randomized DNA-based Crypto-Steganography Algorithm using Generic N-Bits Binary Coding Rule”. Internet of Things and Big Data Technologies for Next Generation Healthcare, Springer International Publishing, 2016 (Accepted).

Table of Contents

Acknowledgment	II
Abstract	III
List of Publications	VII
Table of Contents	IX
List of Figures	XIII
List of Tables	XV
List of Algorithms.....	XVIII
List of Abbreviations	XIX
Chapter 1. Introduction	2
1.1 Overview	2
1.2 Motivation	5
1.2.1 Why Steganography?	5
1.2.2 Why DNA?	5
1.3 Research Objectives	6
1.4 Main Contributions of this Thesis.....	7
1.5 Thesis Organization	9
Chapter 2. Biological Background.....	12
2.1 DNA	12
2.2 Complementary Base Pairing.....	13
2.3 Protein	14
2.4 Mutations.....	16
Chapter 3. DNA Based Steganography Techniques.....	18
3.1 Introduction	18
3.2 DNA Based Steganography Approaches	19
3.2.1 First Approach: Insertion Based Algorithms	20
3.2.2 Second Approach: Substitution Based Algorithms	22
3.2.3 Third Approach: Complementary Rules Based Algorithms	23
3.2.4 Fourth Approach: Combined Approaches Based Algorithms	24
3.3 Issues, Controversies, Problems.....	28
3.3.1 First Approach - Insertion Based Algorithms: M1, M2 and M3	28
3.3.2 Second Approach - Substitution Based Algorithms: M4 and M5	29
3.3.3 Third Approach - Complementary Rules Based Algorithms: M6	30

3.3.4	Fourth Approach: Combined Approaches Based Algorithms	30
3.4	Security Analysis	32
3.5	Comparative Analysis	39
3.5.1	Cryptography Process Factors	40
3.5.2	Steganography Process Criteria.....	44
3.6	Conclusion.....	48
Chapter 4.	DNA Silent Mutations Based Steganography techniques	52
4.1	DNA Based Hybrid Crypto-Steganography Algorithm – Sender Side.....	53
4.1.1	Data Preprocessing Using the Proposed N-Bits BCR ..	53
4.1.2	The First Security Level: The Encryption Layer.....	55
4.1.3	The Second Security Level: The Steganography Layer	58
4.1.3.1	Algorithm 1: DNA LSBBase Substitution Based Steganography ..	58
4.1.3.2	Algorithm 2: The Randomized DNA LSBBase Based Steganography	62
4.2	DNA Based Data Extraction – Receiver Side.....	64
4.2.1	The First Security Level Breaking: The Cipher Data Extraction.....	64
4.2.1.1	Algorithm 1: DNA LSBBase Substitution Based Data Retrieval	64
4.2.1.2	Algorithm 2: The Randomized DNA LSBBase Based Data Retrieval	66
4.2.2	First Security Layer Decoding: Data Decoding	67
4.2.3	Plain Text Formulation	68
4.3	Security Analysis	69
4.3.1	Algorithm 1: DNA LSBBase Substitution Based Steganography	70
4.3.2	Algorithm 2: The Randomized DNA LSBBase Based Steganography	72
4.4	Experimental Results and Discussion.....	74
4.4.1	Dataset and Test Cases	74
4.4.2	Evaluated Parameters.....	76
4.4.3	Results of the Evaluated Parameters.....	76
4.4.4	Comparative Study	78
4.5	Conclusion.....	85

Chapter 5.	High Capacity Conservative Mutations Based Steganography	88
5.1	Conservative Mutations	89
5.2	Conservative Mutations Based Steganography – Sender Side	89
5.2.1	Data Preprocessing Using the Proposed N-Bits BCR ..	90
5.2.2	Lower Security Layer: DNA Based Encryption.....	90
5.2.3	Upper Security Layer: DNA Conservative Mutations Based Steganography	93
5.2.3.1	Why Conservative Mutations?	93
5.2.3.2	The Satisfied Amino Acids	95
5.2.3.3	The Unsatisfied Amino Acids	97
5.2.3.4	Unique unsatisfied amino acids.....	99
5.2.3.5	The Conservative Mutation Based Substitution Steganography	99
5.2.3.6	Example for the Conservative Mutation Based Substitution Rule	101
5.2.4	Illustrative Crypto-Steganography Example - Sender's Work	102
5.3	Conservative Mutations Based Extraction – Receiver Side ..	104
5.3.1	Upper Security Layer: The Decoded Data Extraction	104
5.3.2	Lower Security Layer: The Data Decoding.....	105
5.3.3	Plain Text Formulation	105
5.4	Security Analysis	106
5.4.1	DNA Reference Sequence	107
5.4.2	Binary Coding Rule	107
5.4.3	Playfair Cipher's Matrix	108
5.4.4	The Random Positions of the Hidden Secret Bits	108
5.4.5	Substitution Rule.....	109
5.4.6	The System's Overall Cracking Probability	116
5.5	Experimental Results and Discussion.....	116
5.5.1	Evaluating Parameters for Comparison	117
5.5.2	Test Cases and Data Set.....	117
5.5.3	Results of the Evaluating Parameters	118
5.6	Conclusion.....	121
Chapter 6.	Conclusions and Future Work.....	124
6.1	Cracking Probability	124
6.2	Capacity.....	125

6.3	Numerous Comparison Characteristics.....	126
6.4	Experiments and Results	127
6.4.1	Evaluated Parameters.....	127
6.4.2	Test Cases and Data Set.....	128
6.4.2.1	Comparison Between the Capacity of the PM and Other Approaches	129
6.4.2.2	DNA Sequence's Modification Rate	130
6.4.2.3	Protein's Modification Rate	134
6.4.2.4	The Bits Per Nucleotide Parameter (Bpn)	137
6.4.2.5	The payload parameter	138
6.5	Conclusion.....	138
6.6	Future Work	143
References	145

List of Figures

Figure 2-1 DNA structure	13
Figure 2-2 Amino acids - DNA codons.....	15
Figure 4-1 Hybrid DNA based Crypto-Steganography proposed algorithm phases.....	53
Figure 4-2 The DNA based encryption process flowchart.....	57
Figure 4-3 Hiding sequence bits of M_{bin} and $AMBIG_{bin}$; Row 1 is the real DNA sequence; Row 2 is the least significant carrier base; Row 3 is the LSBase type (Pur. is for Purine, Pyr. is for Pyrimidine); Row 4 is the message bits; Row 5 is the ambiguity bits; Row 6 is the codon's LSBase of the fake sequence; Row 7 is the faked sequence.....	60
Figure 4-4 The steganography phase flowchart	61
Figure 4-5 The hiding process of $CipherM_{bin}$ and $AMBIG_{bin}$ according to random positions (RP) generated using a true random seed number, m_1, \dots, m_6, \dots are the message bits, a_1 and a_2 are the ambiguity bits, S is the original sequence before the steganography while S' is the fake one after it.....	63
Figure 4-6 Hybrid DNA based secret data extraction phases.....	64
Figure 4-7 Cipher data LSBase extraction phase	66
Figure 4-8 The DNA based data decryption and plain text extraction processes flowchart.....	68
Figure 4-9 Modification rate (MR) of the proposed LSBase substitution algorithm against the main substitution one	81
Figure 5-1 The 3 phases of the conservative mutations based crypto- steganography algorithm	90

Figure 5-2 The satisfied amino acids/codons	96
Figure 5-3 Example for the possible substitutions for each unsatisfied amino acid/codons	97
Figure 5-4 DNA conservative mutation based secret data extraction phases	104
Figure 5-5 The possible 6 existing substitution options for Ala's codons if the secret message.....	110
Figure 6-1 The length of the generated sequence after the steganography that is applied by PM VS M1, M2, M3, M4 & M5	130
Figure 6-2 DNA Sequence's Modification Rate	132
Figure 6-3 Protein's Modification Rate.....	134

List of Tables

Table 3-1 An example of the substitution lookup table.....	22
Table 3-2 The six legal complementary rules	26
Table 3-3 Algorithms Comparison: Cryptography Process Factors	42
Table 3-4 Algorithms Comparison: Steganography Process Criteria	47
Table 4-1 4-Bits BCR example to convert a binary message to DNA format	55
Table 4-2 Conversion rule to convert ambiguity numbers (0 to 3) to binary	59
Table 4-3 The eight DNA reference sequences used in the testing phase ...	75
Table 4-4 The evaluated parameters definition	76
Table 4-5 The results obtained to hide 20K bytes of secret data.....	78
Table 4-6 The results obtained using the proposed schemes based on LSBase to hide different sizes of secret data within NW_007906008 real reference sequence	79
Table 4-7 The results obtained using the main substitution algorithm [43] to hide different sizes of secret data within NW_007906008 real reference sequence	80
Table 4-8 A Comparison between the proposed schemes (1 and 2), the main substitution method and the original LSBase method	84
Table 5-1 The Substitution table for the steganography process	101
Table 5-2 The used 4*4 matrix of codons by the DNA based Playfair cipher	102
Table 5-3 Number of the possible substitution rules/Amino acid[Satisfied/Substituted]	113
Table 5-4 Possible substitutions for each amino acid for conservative mutation.....	115