

A Novel Data Hiding Scheme Based on DNA Coding and Module-N Operation

Shun Zhang and Tiegang Gao

*College of Software, Nankai University, PRC, Weijin Road 94#, Nankai District,
Tianjin, 300071, China.
shentengvip@gmail.com*

Abstract

A novel data hiding scheme is proposed based on encoding and substitution. Both the cover media and the secure information are encoded to achieve the substitution. Simulating the genetic codon table in molecular biology, a random grouping algorithm is proposed to generate the pseudo genetic codon table, which regulates the substitution. The cover media is encoded into a DNA codon like sequence. Then, the codons of the cover media in the form of DNA codon sequence are substituted with other codons in the same group according to the pseudo genetic codon table. The secure information in the arbitrary N-nary form decides which codon is selected during every substitution. Due to the random construction of pseudo genetic codon table, encryption is incorporated in the data hiding process. Analysis and experiments have tested the validity and efficiency of the scheme.

Keywords: *Secure communication, Data hiding, DNA coding, Module-N operation*

1. Introduction

Data hiding has wide applications in a variety of areas, such as security communication, watermarking, authentication and so on. Therefore, it has attracted many attentions from both scholars and engineers. Many data hiding schemes have been proposed. For example, LSB (Least Significant Bit) based [1-4], compression based [5], difference expansion based [6, 7] and histogram modification based [8], and so on. Data is embedded just by substituting the least significant bits of pixels in the LSB based data hiding scheme. The compression based data hiding schemes compress the cover media for the embedding room, thus the characters of cover media limit the payload and quality. The difference expansion based schemes hide data by extending the difference between two adjacent pixels [6, 7]. The distortion caused by data hiding increases quickly when the hiding capacity increases. Schemes based on histogram modification hide data by shifting the histograms and it cause less distortion to the cover media. However, the obvious drawback is the limited payload [8].

DNA computing quickly aroused researchers' attention since its first proposal by Adleman [9] in 1994 for solving the Hamiltonian path problem. Due to its good characteristics in parallelism computation, massive storage, and low energy consumption, DNA computing has been widely studied by many scholars in variety of areas such as biology, information science, and mathematics, and so on. In the following years, some NP-complete problems were solved with DNA computing [10-12], which demonstrated its superiority.

Data hiding in DNA utilizes the features of DNA and its translation laws. The earliest secret communication scheme utilizing DNA is proposed by Clelland et al. [13] in Nature. In his scheme, secure data is encoded into DNA like pieces. Then, these pieces are mixed with the ordinary DNA segments to build a pseudo DNA sequence. Finally, these newly built DNA-like sequences are recorded and transmitted to the receiver utilizing microdot

technique. Many data hiding schemes based on DNA has been proposed since then [14-21]. Some schemes just utilize coding way of DNA [16, 17, 22], while others take a full consideration of biology implements [14, 15, 19-21, 23]. The secret information is encode into DNA like sequence and then inserted into existing real DNA sequence in [22]. The length of the DNA sequence is changed after dada hiding, which may reveal the existence of hidden information. Therefore, most schemes keep the length of DNA sequence unchanged [14, 16, 17, 19-21, 23]. Secure data is hidden into DNA sequences by transforming DNA sequences into integers in [17]. Histogram modification is utilized for the data hiding. Data is hidden utilizing the complimentary rule of DNA nucleosides, which means information is embedded in the DNA nucleosides in [16, 22]. However, in [14, 21], data is hidden utilizing the redundancy in the codons translation to amino acids, which means information is embedded in the codons.

All the above-mentioned schemes inspire us to design a general data hiding scheme based on DNA coding and codon substitutions. Therefore, a novel data hiding scheme combining traditional data hiding with DNA encoding is proposed in this paper. Secure information is hidden in to the cover media by the substitution of encoded pseudo codon sequences. There are three highlights in the proposed scheme. The first one is the random construction of pseudo genetic codon table simulating the genetic codon table. The second one is the construction of arbitrary N-nary system according to the pseudo genetic codon table. The last one is the encoding of the cover media and the secure information according to specific rules. The rest of the paper is organized as follows. The preliminaries of the main scheme are proposed in Section 2. Section 3 presents the data hiding algorithm. Section 4 gives the experiments, and Section 5 draws the conclusions.

2. Preliminaries

2.1 Central Dogma and the Construction of Pseudo Genetic Codon Table

The central dogma [24] regulates the transfer of genetic information. The DNA sequence is firstly translated into mRNA. Then, every three nucleosides in the mRNA sequence compose one codon that is mapped to amino acids. Different amino acids are selected according the genetic codon table (Table 1) to construct the proteins. There are 64 codons while there are only 20 amino acids in the genetic codon table. That is to say, there must be several codons mapped to one same amino acid. In fact, there is one, two, four, or six codons mapped to one amino acid, as can be seen in Table 1. Similar to Table 1, other grouping and mapping tables can be randomly generated. Just take Table 2 as an example, denoted by pseudo genetic codon table. The pseudo genetic codon table also maps 64 'codons' to 20 'amino acids'. Table 2 is quite similar to Table 1, as they both map a set of codons with size 64 to a set of codons with size 20. The pseudo genetic codon table can be of various styles using the randomizing technique. Suppose there are S kinds of codes in the set C , and they are to be partitioned into N groups, where $N \leq S$. Here is the algorithm for randomly generating such a table as Table 2: Generate a random string Y containing N integer numbers in a range of $[t_1, t_2]$, where $0 \leq t_1 < t_2 \leq S$, and obviously $S = \sum_{i=1}^N Y(i)$. Select $Y(i), 1 \leq i \leq N$ codes from set C as the i th group. Finally, rearrange these groups in the genetic codon table form. The generated table will be different from these two tables as S and N change.

As in Table 2, the corresponding elements in set C are:

{000,001,002,003,100,101,102,103,200,201,202,203,300,301,302,303,010,011,012,013,110,111,112,113,210,211,212,213,310,311,312,313,020,021,022,023,120,121,122,123,220,221,222,223,320,321,322,323,030,031,032,033,130,131,132,133,230,231,232,233,330,331,332,333}; and the size $S=64$. They are catalogued into $N=20$ groups, which are {A,B,C,E,F,G,H,I,J,K,L,M,N,O,P,Q,R,S,T} ; and the random string $Y=[4,4,4,5,3,6,5,6,5,2,2,2,2,2,2,2,2,2,2,2]$.

Table 1. Standard Genetic Codon Table

1st base	2nd base				3rd base
	U	C	A	G	
U	UUU (Phe/F)	UCU (Ser/S)	UAU (Tyr/Y)	UGU (Cys/C)	U
	UUC	UCC	UAC	UGC	C
	UUA	UCA	UAA Stop	UGA Stop	A
	UUG	UCG	UAG	UGG (Trp/W)	G
C	CUU (Leu/L)	CCU (Pro/P)	CAU (His/H)	CGU	U
	CUC	CCC	CAC	CGC	C
	CUA	CCA	CAA (Gln/Q)	CGA	A
	CUG	CCG	CAG	CGG	G
A	AUU (Ile/I)	ACU (Thr/T)	AAU (Asn/N)	AGU (Ser/S)	U
	AUC	ACC	AAC	AGC	C
	AUA	ACA	AAA	AGA	A
	AUG (Met/M)	ACG	AAG (Lys/K)	AGG (Arg/R)	G
G	GUU	GCU (Ala/A)	GAU (Asp/D)	GGU	U
	GUC	GCC	GAC	GGC	C
	GUA	GCA	GAA (Glu/E)	GGA	A
	GUG	GCG	GAG	GGG	G

Table 2. Pseudo Genetic Codon Table

1st base	2nd base				3rd base				
	0	1	2	3					
0	000	A	010	B	020	C	030	D	0
	001		011		021		031		1
	002		012		022		032		2
	003		013		023		033		3
1	100	H	110	G	120	F	130	E	0
	101		111		121		131		1
	102		112		122		132		2
	103		113		123		133		3
2	200	J	210	I	220	K	230	L	0
	201		211		221		231		1
	202		212		222		232		2
	203		213		223		233		3
3	300	N	310	R	320	O	330	P	0
	301		311		321		331		1
	302		312		322		332		2
	303		313		323		333		3

2.2 DNA Encoding for the Cover Media

As is known, there are four kinds of nucleosides in DNA sequence or RNA sequence. Moreover, two binary bits can represent four states. Therefore, binary streams can be encoded into DNA or RNA sequences. Images are the most frequently used in daily life. Just take the image as the example of cover media in describing the encoding procedure. Suppose there is a two-dimensional cover image I with size $M \times N$. The detailed encoded steps are presented as follows: Firstly, scan image I to get one-dimensional sequence S with size $(1, M \times N)$, the scanning method can be the simple ‘from left to right and from top to bottom’ method, or other ways such as the “zigzag” way or “invers S ” way. Secondly, transform every grey value $S(i), i = 1, 2, \dots, M \times N$ into binary sequence $S(i, 1:8)$. Finally, encode the last $k, (1 \leq k \leq 8)$ bits $S(i, (8-k):8)$ of every $S(i, 1:8)$ into the pseudo DNA sequence DNA according to any of the mapping ways as presented in Table 3. Every three adjacent nucleosides construct one codon. Note that, all the k -bit segments are concatenated to form a long bit sequence, which is then transformed to the sequence of codons in practice.

Table 3. Different Mapping Ways in DNA Coding

Codes	1	2	3	4	5	6	7	8
00	A	A	C	C	G	G	T	T
01	C	G	A	T	A	T	C	G
10	G	C	T	A	T	A	G	C
11	T	T	G	G	C	C	A	A

2.3 Arbitrary N-nary Transformation for the Secure Information

From the perspective of computer science, all information can be presented with binary streams through sampling and coding. Therefore, random binary streams are adopted as the secret data. As is known, binary streams can be transformed into decimal, octal, and hexadecimal form. Therefore, it can also be transformed into N-nary form, where N is arbitrary. Denote such a N-nary system as arbitrary N-nary system. Suppose that there is an arbitrary N-nary system with bases $P = \{P_i, i = 1, 2, \dots, m\}$. Obviously, the numerical codes for different bases are numbers in ranges of $[0 \square P_i - 1]$ respectively. Binary streams can be transformed into the arbitrary N-nary form with a simple modulation operation. Suppose there is a binary stream B . Firstly, transform B into decimal form D . Then, calculate the i th numerical code $M(i)$ with $M(i) = \text{mod}(B / \prod_{k=1}^{i-1} P_k, P_i)$, where $i \geq 1$.

3. Data Hiding and Extraction Scheme

Based on the encoding and transformation strategies proposed in the above section, the main process of the data hiding scheme is substitution. Before the substitution, the cover media is encoded into a DNA like sequence (pseudo DNA sequence). Then, it can be rearranged into a codon like sequence (pseudo codon sequence). After that, a pseudo genetic codon table like Table 2 is constructed for mapping the pseudo codon sequence into different groups. Every codon in the pseudo codon sequence belongs to one group in the pseudo codon table. The amount of codons in one group is random according to the randomly constructed pseudo genetic codon table. Therefore, with the pseudo genetic codon sequence, an arbitrary N-nary system is achieved like Table 4. In Table 4, the original pseudo codon sequence is 000 032 202 131 332 101. According to pseudo genetic codon table (Table 2), their corresponding groups of the codons are 000,001,002,003 , 030,031,032,033,130 , 202,203 , 131,132,133 , 332,333 , 100,101,102,200,201 respectively. The bases of the constructed arbitrary N-nary system are $P = \{4,5,2,3,2,6\}$. Secure data to be embedded is transformed into N-nary form by the module-N operation according to the constructed arbitrary N-nary system. Finally, the pseudo codons in the pseudo codon sequence are substituted with codons in the same group according to the N-nary formatted information to achieve the data hiding. The substituted pseudo codon sequence with information embedded is decoded to its original form. By selecting suitable pseudo genetic codon table, slight changes will be introduced to the cover media, and high embedding rate can be achieved. The embedded information can never be extracted without the randomly generated pseudo genetic codon table, which is generated by random strings.

Table 4. Arbitrary N-nary System and Codon Substitution

M(j)	A	D	J	E	T	H
0	000	030	202	131	332	100
1	001	031	203	132	333	101
2	002	032		133		102
3	003	033				103
4		130				200
5						201

Detailed data hiding algorithm and data extraction algorithm are presented in the followings.

3.1 Data Hiding Algorithm

- Encode the cover media into pseudo codon sequence C with the method proposed in section 2.2;
- Generate the pseudo genetic codon table with the method proposed in section 2.1;
- Build an arbitrary N-nary system with the pseudo DNA codon sequence according to the pseudo genetic codon table;
- Transform the secret information S into binary sequence B , then transform B into decimal integer D , finally transform D into stream M in arbitrary N-nary format through the method proposed in section 2.3;
- Substitute the pseudo codons with the $M(j)$ th codon in the same group;
- Translate the substituted pseudo codon sequence into original form.

Parameters for encoding the cover media and generating the pseudo genetic codon table are encrypted as keys for information extraction, denoted by $key1$ and $key2$ respectively.

3.2 Data Extraction Algorithm

- Transform the received cover media into pseudo codon sequence with $key1$;
- Generate the pseudo genetic codon table with $key2$;
- Build the arbitrary N-nary system according to the pseudo codon sequence and pseudo genetic codon table;
- Read the N-nary sequence M according to the corresponding position of current codon in the arbitrary N-nary system table (Table 4);
- Transform M into binary form and decode it into original information.

To make the scheme clearly understood. Some key parameters are described as follows. The encoded pseudo codon sequence is $C=\{000,032,202,131,332,101\}$, and the secure information is $B=\{1010010101\}$. An arbitrary N-nary system is built like Table 4, and its bases are $P=\{4,5,2,3,2,6\}$. The decimal form of B is $D=661$. The arbitrary N-nary stream $M=\{1,4,0,1,0,1\}$, and detailed process to get M : $M(H)=\text{mod}(297/6)$, $M(T)=\text{mod}(297/6/2)$, $M(E)=\text{mod}(297/6/2/3)$, $M(J)=\text{mod}(297/6/2/3/2)$, $M(D)=\text{mod}(297/6/2/3/2/5)$, $M(A)=\text{mod}(297/6/2/3/2/5/4)$. Therefore, the substituted pseudo codon sequence $C'=\{001,130,202,132,332,101\}$.

4. Experiments and Analysis

To testify the efficiency and validity of the proposed scheme, images (with size 512×512) from USC-SIPI image database, Miscellaneous gray level images and UCID data base are selected as the cover media for experiments. Random binary stream are embedded as the secret data. All experiments are performed on the MATLAB 2012a platform running on a personal computer with a CPU of AMD Phenom (tm) II X4 810 @ 2.6GHz, memory of 4 GB, and operating system of Windows 7 x64 Ultimate Edition.

Two typical images ‘Airplane’ and ‘texture’ after data hiding are presented in Fig. 1. The parameter k is often selected. The last 6 LSB and 2 LSB are selected as the cover media for data hiding. The embedding rates of sub-image (a) and (b) in Fig.1 are 2.0310 bpp (bit per pixel) and 0.5883 bpp respectively, while the embedding rate of image (c) and (d) are 1.7992 bpp and 0.5938 bpp respectively.

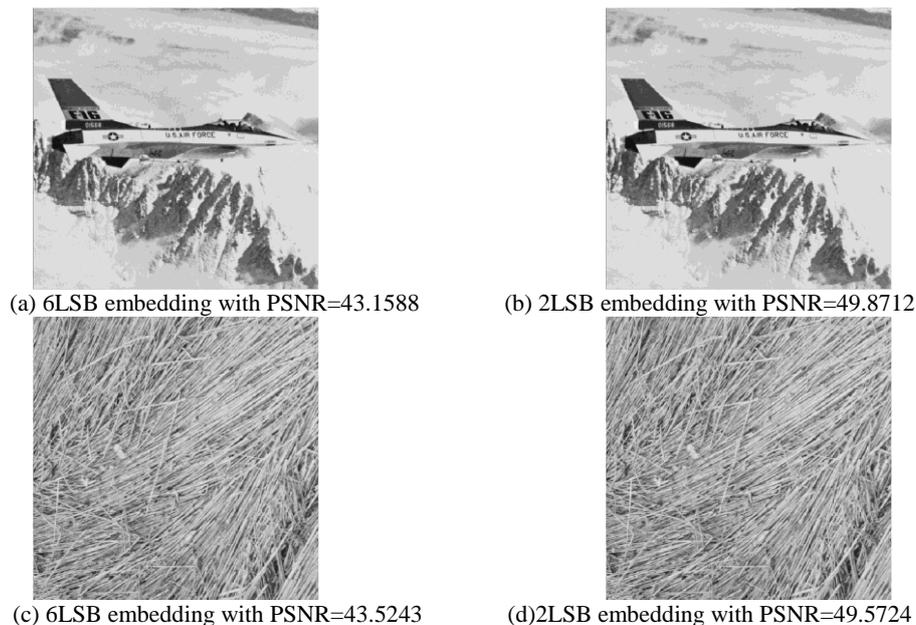


Figure 1. Image ‘Airplane’ and ‘Texture’ after Embedding

The payload (embedding capacity) varies as the cover media changes, because different cover media will be encoded into different pseudo codon sequences. Besides, the payload varies when different pseudo genetic codon tables are generated. When the pseudo genetic table is Table 2, the embedding rate and corresponding PSNR of different images are presented in Table 5. The last six bits of every pixel in the image are encoded for the information hiding.

Table 5. Results of Last Six Bits Embedding According to Table 2

Image	Airplane	Baboon	Boat	Lena	Peppers	Tiffany
ER(bpp)	2.03	1.73	1.97	1.81	1.87	1.70
PSNR(dB)	43.16	43.47	43.70	44.03	43.78	43.65

Compared with those LSB based data hiding schemes [1-4], the proposed scheme offers more security through the randomly generation of pseudo genetic codon table. One can never know the arbitrary N-nary system without the pseudo genetic codon table that is generated by the random grouping of codons. The random substitutions of pseudo codons cause random changes to the cover media, which means that it resist attacks base on statically analysis, such as χ^2 analysis. Besides, as can be seen in Table 4, information is embedded without changing the original pseudo codons sometimes. Therefore, the proposed scheme cause less distortion. Experimental results in Table 6 and Table 7 also proved it. The last six bits of every pixel are selected for the data hiding in Table 6, while in Table 7, the last two bits of every pixel are selected for the data hiding.

Table 6. Comparisons between LSB Algorithm and Scheme by Us

Image	Airplane	Baboon	Boat	Lena	Peppers	Tiffany	
ER(bpp)	1.50	1.50	1.50	1.50	1.50	1.50	
PSNR(dB)	LSB	45.42	45.40	45.26	45.39	45.41	45.41
	Ours	46.93	46.47	46.20	47.09	46.98	46.74

Table 7. Comparisons between LSB Algorithm and Scheme by Us

Image	Airplane	Baboon	Boat	Lena	Peppers	Tiffany	
ER(bpp)	0.3	0.3	0.3	0.3	0.3	0.3	
PSNR(dB)	LSB	52.3878	52.4339	52.2734	52.3576	52.4145	52.4184
	Ours	52.7312	52.9913	52.7304	52.7862	53.0851	53.1859

Higher embedding capacity and lower distortion has been achieved compared with those traditional data hiding schemes such as difference expansion based methods [6, 7], histogram based methods [8], and compression based methods [5].

The DNA data hiding schemes based on complementary rule [16, 22] hide data through substitution of the nucleosides with their complementary nucleosides. Reference DNA sequence is needed when data extraction. However, scheme proposed by us achieves blind extraction that is very important in data hiding. Besides, in [16], the lookup table is required in data extraction, and thus additional information needs to be transmitted. Moreover, the look up table makes the scheme unsafe. The hidden information will be easily extracted if the look up table is intercepted.

Embedding efficiency is defined as the number of information bits embedded as one change in the cover media occurs [25]. Twenty images are selected from USC-SIPI image database to calculate the embedding efficiency. The embedding efficiency of the proposed scheme is 3.6134 while the PSNR between original images and images with information embedded is 43.2164 on average.

Suppose the genetic codon table is just the real one (Table 2). The core of proposed scheme is partly coincidental to DNA data hiding method proposed by Tai *et al.* [14]. In fact, the proposed scheme can be used for any digital cover media as long as it can be encoded, DNA sequence is of course one of them. In his scheme, marks are embedded into mRNA codons through a substitution of mRNA codons that can be mapped to the same amino acids for the authentication of plant variety rights. Scheme by Tai *et al.* [14] depends on the real genetic codon table that is fixed and public. Besides, original mRNA codon sequence is required when extract the marks. However, scheme proposed in this paper depends on the pseudo genetic codon table that is generated randomly with a one-time pad key.

5. Conclusions

A cover media encoding scheme simulating the central dogma in biology and an arbitrary N-nary transformation of additional information are proposed for the design of secure data hiding. Both the cover media and the information to be embedded are encoded, and then substitution of elements in cover media is imposed to represent the data to be embedded. The high embedding capacity and low distortion are analyzed in theory and demonstrated through experiments. Because the embedding of data relies on the pseudo genetic codon table that is generated through random segmentation, the encryption is realized along with data embedding. Besides, not only images but also any other digital media can be encoded as the cover media in the proposed data hiding scheme.

Acknowledgements

The work described in this paper was supported by the Key Program of National Science Fund of Tianjin, China (Grant NO. 11JCZDJC16000).

References

- [1] R G. Van Schyndel, A Z. Tirkel and C F. Osborne, "A digital watermark", in Image Processing, Proceedings, IEEE International Conference, (1994), pp. 86-90.
- [2] C-K. Chan and L-M. Cheng, "Hiding data in images by simple LSB substitution", Pattern recognition, vol. 37, no. 3, (2004), pp. 469-474.
- [3] C-H. Yang, C-Y. Weng and S-J. Wang, "Adaptive data hiding in edge areas of images with spatial LSB domain systems", Information Forensics and Security, IEEE Transactions on, vol. 3, no. 3, (2008), pp. 488-497.
- [4] X. Liao, Q-y. Wen and J. Zhang, "A steganographic method for digital images with four-pixel differencing and modified LSB substitution", Journal of Visual Communication and Image Representation, vol. 22, no. 1, (2011), pp. 1-8.
- [5] M U. Celik, G. Sharma and A M. Tekalp, "Lossless generalized-LSB data embedding", Image Processing, IEEE Transactions on, vol. 14, no. 2, (2005), pp. 253-266.
- [6] J. Tian, "Reversible data embedding using a difference expansion", Circuits and Systems for Video Technology, IEEE Transactions on, vol. 13, no. 8, (2003), pp. 890-896.
- [7] A M. Alattar, "Reversible watermark using the difference expansion of a generalized integer transform", Image Processing, IEEE Transactions on, vol. 13, no. 8, (2004), pp. 1147-1156.
- [8] Z. Ni, Y-Q. Shi and N. Ansari, "Reversible data hiding", Circuits and Systems for Video Technology, IEEE Transactions on, vol. 16, no. 3, (2006), pp. 354-362.
- [9] L M. Adleman, "Molecular computation of solutions to combinatorial problems", Science, New Series, vol. 266, no. 5187, (1994), pp. 1021-1024.
- [10] R J. Lipton, "DNA solution of hard computational problems", Science, vol. 268, no. 5210, (1995), pp. 542-545.
- [11] Q. Ouyang, P D. Kaplan and S. Liu, "DNA solution of the maximal clique problem", Science, vol. 278, no. 5337, (1997), pp. 446-449.
- [12] Q. Liu, L. Wang and A G. Frutos, "DNA computing on surfaces", Nature, vol. 403, no. 6766, (2000), pp. 175-179.
- [13] C T. Clelland, V. Risca and C. Bancroft, "Hiding messages in DNA microdots", Nature, vol. 399, no. 6736, (1999), pp. 533-534.

- [14] W.-L. Tai, C. C. N. Wang and P. C. Y. Sheu, "Data Hiding in DNA for Authentication of Plant Variety Rights", *Journal of Electronics*, vol. 11, no. 1, (2013), pp. 38-43.
- [15] D. Haughton and F. Balado, "BioCode: Two biologically compatible Algorithms for embedding data in non-coding and coding regions of DNA", *BMC bioinformatics*, vol. 14, no. 121, (2013), pp. 1-16.
- [16] J-S. Taur, H-Y. Lin and H-L. Lee, "Data hiding in DNA sequences based on table lookup substitution", *International Journal of Innovative Computing Information and Control*, vol. 8, no. 10 A, (2012), pp. 6585-6598.
- [17] Y-H. Huang, C-C. Chang and C-Y. Wu, "A DNA-based data hiding technique with low modification rates", *Multimedia Tools and Applications*, vol. 70, no. 3, (2014), pp. 1439-1451.
- [18] D. Tulpan, C. Regoui and G. Durand, "HyDEn: A Hybrid Steganocryptographic Approach for Data Encryption Using Randomized Error-Correcting DNA Codes", *BioMed research international*, vol. 2013, (2013), pp. 1-11.
- [19] D. Heider and A. Barnekow, "DNA-based watermarks using the DNA-Crypt algorithm", *BMC bioinformatics*, vol. 8, no. 1, (2007), pp. 176.
- [20] P.C. Wong, K-k. Wong and H. Foote, "Organic data memory using the DNA approach", *Communications of the ACM*, vol. 46, no. 1, (2003), pp. 95-98.
- [21] B. Shimanovsky, J. Feng and M. Potkonjak, "Hiding data in DNA," in *Information Hiding*, (2003), pp. 373-386.
- [22] H. Shiu, K-L. Ng and J-F. Fang, "Data hiding methods based upon DNA sequences", *Information Sciences*, vol. 180, no. 11, (2010), pp. 2196-2208.
- [23] D. Tulpan, C. Regoui and G. Durand, "HyDEn: A Hybrid Steganocryptographic Approach for Data Encryption Using Randomized Error-Correcting DNA Codes", *BioMed research international*, vol. 2013, (2013).
- [24] F. Crick, "Central dogma of molecular biology", *Nature*, vol. 227, no. 5258, (1970), pp. 561-563.
- [25] J. Fridrich and D. Soukal, "Matrix embedding for large payloads," in *Electronic Imaging 2006*, (2006), pp. 60721W-60721W-12.

Authors



Shun Zhang, he born in 1986, is a Ph.D. candidate in in College of Software, Nankai University, China. His current research interests include steganography, data hiding, and multimedia security.



Tiegang Gao, he is a Prof. in College of Software, Nankai University, China since 2006. His current research interests include information security, multimedia information processing and software engineering.